# The Object at Hand: Automated Editing for Mixed Reality Video Guidance from Hand-Object Interactions

Yao Lu*
University of Bristol

Walterio W. Mayol-Cuevas†
University of Bristol

## ABSTRACT

In this paper, we concern with the problem of how to automatically extract the steps that compose real-life hand activities. This is a key competence towards processing, monitoring and providing video guidance in Mixed Reality systems. We use egocentric vision to observe hand-object interactions in real-world tasks and automatically decompose a video into its constituent steps. Our approach combines hand-object interaction (HOI) detection, object similarity measurement and a finite state machine (FSM) representation to automatically edit videos into steps. We use a combination of Convolutional Neural Networks (CNNs) and the FSM to discover, edit cuts and merge segments while observing real hand activities. We evaluate quantitatively and qualitatively our algorithm on two datasets: the GTEA [18], and a new dataset we introduce for Chinese Tea making. Results show our method is able to segment hand-object interaction videos into key step segments with high levels of precision.

**Index Terms:** Human-centered computing—Human computer interaction (HCI)—Interaction paradigms; Applied computing—Education—Computer-assisted instruction

## 1 INTRODUCTION

Guidance is one of the core target applications for Augmented Reality (AR), specifically in industry or to support daily living. AR systems have often been proposed as enablers to perform maintenance and repair tasks and several systems have been devised for this task [2, 10, 24].

However, content authoring for real-world tasks is one of the principal challenges in the extended Mixed Reality (MR) field, and one that limits the adoption of MR technology in general [16, 30]. Overall, automating content authoring for MR systems has received relatively little attention in contrast to the body of work that explores how such systems could be used. Commonly, content creation assumes pre-definition of object and scene model assets. This makes content authoring expensive, inefficient and restricted to the manual workflow and content authoring tools available. It is also common to impose strict sequences of pre-defined steps on activities. However, it is precisely real-world tasks such as those in daily living, that have limited pre-existing assets and are largely unscripted that can be helped the most by systems that support memory loss, DIY tutoring, maintenance and assembling tasks. In this paper we argue that it is thus better to work-backwards and observe the world as-is, extracting the task steps and relevant objects that real people use in real tasks and in the wild.

On the other hand, First Person View (FPV) content is produced with increased ease due to the fast development of wearable devices and video-based social media. An egocentric system follows the user's viewpoint, helping to record or deliver data for tutoring of a

---

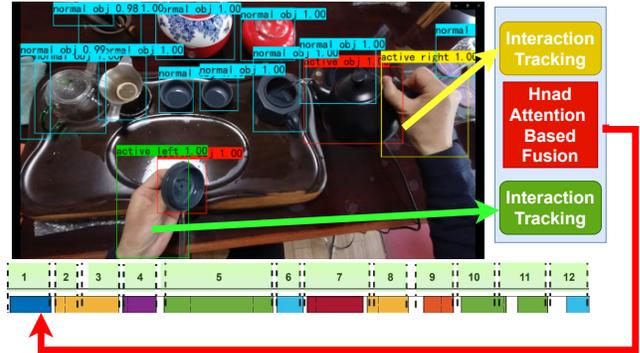*e-mail: yl1220@bristol.ac.uk

†e-mail: wmayol@cs.bris.ac.uk

Figure 1: Our hand-object interaction (HOI) system is able to identify objects being manipulated (two red boxes), each of the hands (yellow and green boxes), as well as ignore other objects in the scene (cyan). Using this information we can segment a task into its constituent steps (time-graph below image) for guidance model building or task monitoring.

task to the same or different user. Overall, humans are well-tuned to imitate from watching and learning from imitation, a survival competence that is present in other animals too [14]. In terms of recent video guidance systems for MR, the system in [16] showcased that using egocentric video guidance is easier to produce and more accessible for users to follow instructions while freeing their hands. In [20], results show a video-based AR system able to show guidance and monitor hands but from pre-defined hand poses and locations. In this work, we expand from prior MR content generation work by focusing on the critical aspect of automated task decomposition via hand-object interaction detection.

Our system (shown in figure 1) is not an end-to-end closed box in the fashion of other Deep Learning works. Here, we do use CNNs for the relevant aspects of the task decomposition such as hand and object detection in the wild but combine these with algorithms that help explainability, such as a finite state machine for hand-object interaction control.

Our proposal in this paper detects hand-object interactions from egocentric video of real-world tasks and segments the video into task segments and hand-object events. This is the first crucial step towards the long term objective of assisting in task guidance. Our contributions are threefold: 1. We propose and implement a method that can identify hands, task-relevant objects and hand-object interactions from egocentric videos in the wild. 2. A method to automatically decompose the input video into its constituent activity steps. 3. Evaluation of our method on non-scripted egocentric datasets quantitatively and qualitatively, including the introduction of a new dataset of Chinese Tea making with expert annotations.

The rest of our paper is organized as follows: After reviewing the related work in section 2, we describe our content authoring paradigm in section 3. We then evaluate our results in section 4 and present the conclusions in section 5.

## 2 RELATED WORK

The key contribution of this work is automatic video content authoring with hand-object interaction cues. In this section, we mainly look at work related to the field of unsupervised video content authoring and egocentric hand-object interaction detection.

### 2.1 Video Content Authoring and Guidance

Automating the edit of a video as a tutorial can be traced to systems like Kang and Ikeuchi's [13] that classified the hand-object interaction into five phases: approaching, grasp phase, manipulation phase, place object and depart. By analysing fingertip polygon and hand movement speed, breakpoints can be found within an interaction process. Their work is very inspiring, and to our knowledge, it is the first work that utilising hand as a cue to segment an interaction. The work from Mayol and Murray [21] took one step forward with the objects being interacted with hands extracted by an attention filter that used an in-situ learned skin colour distribution. This allowed to extract a pictorial summary of the interaction. The work from Michihiko et al. [8] overlays 2D video guidance information in a 3D viewpoint. Overall, there is growing interest in video-based guidance in MR since it can be shown that people can be supported by a system that can display step-wise video guidance on a headset and video guidance reduces the hardware and UX requirements of such system compared to fully-fledged positional 6D AR. The work of Nils and Didier [23], proposed another unsupervised workflow video segmentation method. Instead of using high-level features like hand and object interactions, they build a distance function with region descriptors to segment the video of a simplified repair task according to the centre crop of each frame. Other systems have looked at how to combine multiple observations for task decomposition. In, Longfei et. al. consider manipulation "hotspots" to identify critical parts of hand-object interactions from multiple observations on a specific object, a sewing machine. A key first step on all the above systems towards video-based guidance is that the content to be displayed needs to be segmented into coherent steps in advance, something that has been a challenge when considering true non-scripted, multi-object tasks in the wild. On the other hand there have been recent egocentric and non-egocentric deep learning based video segmentation methods such as [6, 12, 15, 17, 28, 29] that showcase good ability to segment videos into actions. These methods are supervised and heavily rely on the training data, more details can be found in the survey [1].

### 2.2 Interaction Relevant Object Detection

Detecting relevant task objects for interaction has received substantial interests in the computer vision community. According to how the interaction is defined, this problem is solved differently. For instance, the early work of [21] considers as the task-relevant object appearing at the centre of the image frame over a few seconds. While in Teesid et al. [16], they hypothesize that the task-relevant object locates within the gaze region that can be predicted with a head-mounted IMU (Inertial Measurement Unit). However, the way of finding task-relevant objects in [16] is closer to scene recognition. With promising performance achieved by CNN-based object detectors such as [7] [26] [25], it is possible to train an end-to-end model for object-hand interaction detection. Like the work of Dandan et al. [27], they directly label the hand status and active object for a Faster-RCNN [26] based network for detection. This is also what we argue and expand in this work. The interaction/task-relevant object is highly related to the human hand, especially in First Person View. We can assume the hands in view with the right pose relative to the camera mostly belong to the camera wearer. Moreover, the objects interacting with these hands are task-relevant. In our work, we follow a similar paradigm as Shan et al. [27] for task-relevant object detection. But instead of considering all possible hand-object
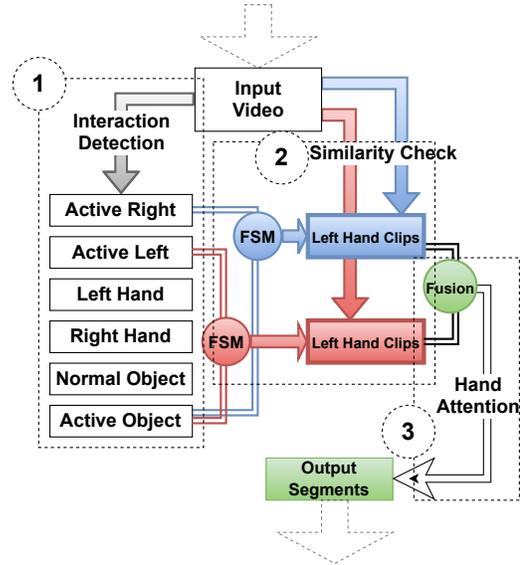


Figure 2: A schematic of our automatic content authoring system. Our system runs in three main stages: (1) Hand-object interaction (HOI) extraction, hand based clips segmentation and segments fusion. With an input video, HOI information is extracted by a Faster-RCNN based detector. (2) By controlling the hand status with FSM (finite state machine) separately, two streams of segmentation can be obtained. (3) Finally, with the hand attention prediction, two streams containing the same step can be fused into one (See Fig 8).

cases, we focus on the first-person view and train the network with our hand object interaction dataset. More details in section 3.2.

## 3 CONTENT AUTHORING SYSTEM

Figure 2 shows a schematic of our system. The workflow of our system is straightforward and runs in three stages: HOI (hand-object information) extraction, hand-based clips segmentation and segments fusion. Given that an input video captured from an expert performs the task in a proper sequence, a Faster-RCNN [26] detector is applied to extract the HOI information from each frame. The detected objects and hands are fed into our FSM (Finite State Machine) for hand status control. We can obtain a two-stream initial video segmentation by cutting from the point when the hand status changes from two hands. To prevent over-segmentation of videos, we use a CNN based similarity network to measure the similarity between two adjacent clips. As an outcome, two streams of segmentation (from the left and right hands) are merged into one according to their IOSA (intersect over the smaller area) over time and our hand attention prediction. The result is an input video segmented into its constituent steps.

### 3.1 Problem Formulation

We denote the full input video as $V = \{I_1, ..., I_T\}$ which has $T$ frames, for each frame $I_t$, we perform our hand-object interaction detector and obtain a set of bounding boxes $Box^{I_t} = \{b_{n\_o}, b_{a\_o}, b_{n\_l}, b_{n\_r}, b_{a\_l}, b_{a\_r}\}$ which stand for 'normal object' (an idle object not being interacted with), 'active object' (an object which is being interacted with), 'idle left hand' (left hand which is not interacting with an active object), 'idle right hand', 'active left hand' (object manipulation by left hand), 'active right hand' (object manipulated by right hand). Taking a practical approach to defining hand tasks observed by egocentric video, we assume for a given frame there is *at most* 1 left hand, 1 right hand and 2 active objects (i.e. at most one active object per hand). We use a FSM (Finite
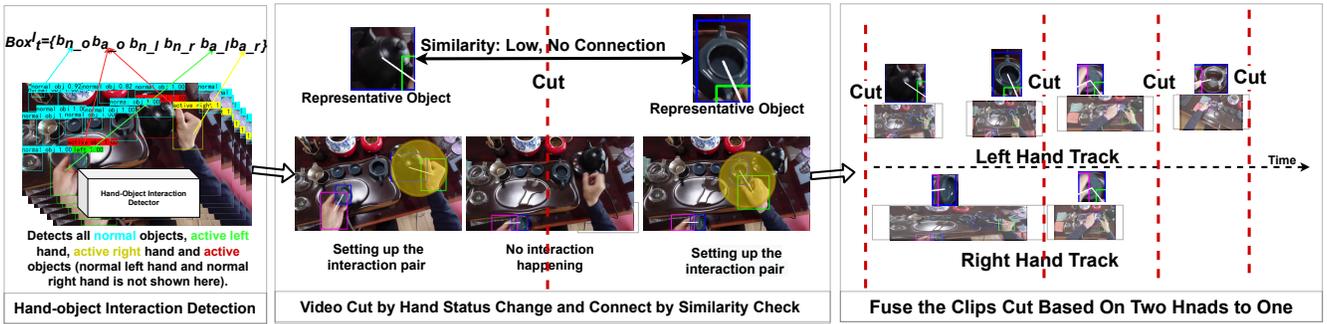
Figure 3: An intuitive illustration of the content authoring process. Left: Hand-Object Interaction information is extracted. The information contains the bounding boxes and classification of all objects and hands. Middle: These extracted information is fed to a finite state machine controlling the hand status (not shown in the picture). Once the hand status changes (e.g., status change from 'interacting with object' to 'idle' or vice versa), we introduce a video segment split. To prevent the over-segmenting of the video, we measure the similarity of two clips based on the automatically obtained object crops and reconnect two adjacent clips if they have high similarity. Right: The video is cut into two streams according to hand and object status. By calculating the IOSA over time, we can unify segments that belong to the same task step. More details on fusion can be found in figure 8

State Machine) to control the status $S_l^t$ for a hand. The hand status changes with different input from HOI detection results. The initial segmentation can be obtained by splitting the video where the hand status changes. The initial segmentation of left hand stream can be represented as: $C_l = \{c_1, ..., c_m\}$, where $m$ is the number of clips. We use a high threshold in HOI detection in each frame to reduce the risk of accumulating the error across the video. However, the video could be over-segmented due to miss detected changes of hand status. As compensation, to reconnect similar clips, we measure the average similarities $\bar{E}$ between the adjacent clips in $C$:

$$\bar{E}_{m-1,m} = \frac{\sum_{i=1}^{x} \sum_{j=1}^{y} s_{ij}}{i \times j} \quad (1)$$

where $x$ and $y$ are the total number of images chosen from two clips, $s_{ij}$ is a CNN based similarity measure between two images inspired by Siamese [22] network. If $\bar{E}_{m-1,m}$ is less than a threshold $T$, we consider the two clips as one. $T$ is found experimentally and based on the performance of all the datasets we use here. The initial segmented clip set of right hand $C_r$ can be obtained in the same way. In the end, according to the IOSA (interaction over the smaller area) over time between the two-stream clip sets and the hand attention prediction. We obtain the final segmentation $C_{fused}$.

## 3.2 Hand-Object Interaction Detection

### 3.2.1 Faster-RCNN Based Detector

Hand-object interaction is the principal cue in the content authoring of our system. Given one frame from an egocentric view, we consider at most one left hand and one right hand. Each hand can handle at most one object. The hand-object interaction is defined when a hand is physically close to an object. We define a hand that is interacting with an object as an 'active hand'-'active object' pair. Otherwise, the hands and objects in the scene without interaction are defined as 'idle hand' and 'normal object'. We detect every possible 'object' (including hands) in the scene. Each 'object' is classified as 'active left/right hand', 'idle left/right hand', 'active object' or 'normal object'. Our hand-object interaction detector is a trained ResNet-backbone [11] Faster-RCNN network [26]. It takes an image as input and outputs bounding boxes of 6 possible classes: 'normal object', 'active object', 'active left hand', 'active right hand', 'idle left hand' and 'idle right hand'. As far as we are aware, there is only one public dataset from [27] that has a similar labelling pattern. However, their label is 'hand' centred, the 'active object' is assigned

to 'hand', and no 'normal object' label is provided. In addition, most of their data are third-person views which are not suitable for our application. We, therefore, had to do our own labelling for the datasets we use in this paper. We are publicly releasing labels, our new dataset and code accompanying this publication.

### 3.2.2 FPV Hand-Object Interaction Dataset

In total, we label $3K$ images with 6 labels as per above. Part of the labelled images are chosen from first-person view based datasets, including EPIC-KITCHENS [3], GTEA [18], and First-Person Hand Action [5]. We believe that in some applications, training on a small, specially targeted dataset could make a detector outperforms those trained on a larger general-purpose dataset. In the data selection stage, we intentionally picked images where each:

1) has hands absent or idle
2) has one hand interaction and the other idle or absent
3) has two hands interaction separately or together
4) has one hands just about to interact
5) has hand-object interaction with cluttered background
6) has hand-object occlusion

Example labelled images are shown in Figure 4. All objects are labelled for each image. We found this very important for fine-tuning the pre-trained RPN (region proposal network) to discover possible objects that have not been seen before. Furthermore, detecting all possible objects is helpful in guidance delivery for step-object compliance checks.

## 3.3 Initial Segmentation

Here we discuss how we utilize the extracted HOI information for hand status identification.

### 3.3.1 Hand Status Control

After obtaining frame-wise detection results $Box^{I_t}$, we build a FSM based hand status model. In the work of Kang and Ikeuchi [13], they consider five phases for hand-object interaction: 'approaching', 'grasp', 'manipulation', 'place down' and 'depart'. This model is object-centred and too idealized to apply in realistic settings. For example, the approach and depart stages are not always clear when the hand is very close to the object without additional information. Also, the boundary between 'grasp' and 'manipulation' is hard to define. Thus, our FSM only considers two different states: ' active' and 'idle' (representing interaction and non-interaction). To have
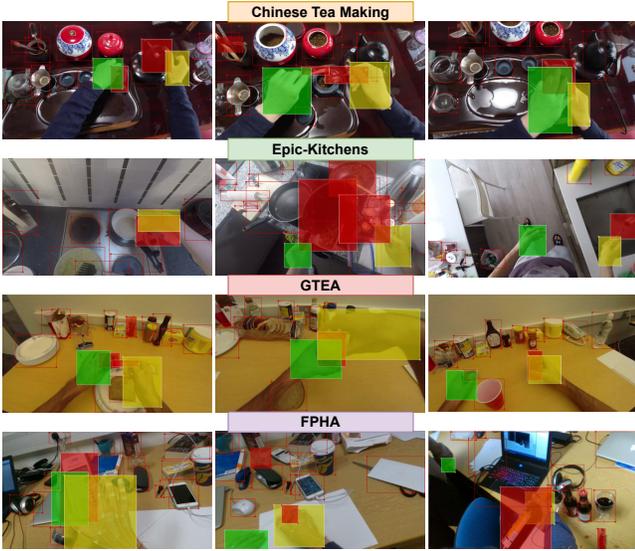
Figure 4: The demo of labelled images from our Chinese Tea Making dataset (row 1), EPIC-Kitchens [3] dataset (row 2), GTEA [18] dataset (row 3) and FPHA [5] dataset (row 4). Filled green boxes represent left hands, filled yellow boxes represent right hands (activeness is not reflected in the figure). Filled red boxes are 'active object' and other red boxes are 'normal object'.

a clear-cut between the two states, we collect and train our hand-object interaction detector with the negative label on 'hand-object approaching' and 'hand-object departing' cases.

In our FSM, considering the left hand in a frame, if the IOU of bounding boxes between any 'active left hand' and 'active object' is greater than zero, we take the 'score' of this frame as '1' (reflected in Figure 5a). A sliding window of length $n$ is applied to sum up all 'scores' within it. If the summation from frame $t - n$ to $t$ is greater than a threshold $T$, at the frame $t$, hand status is 'active' (as illustrated in Figure 5b). Parameters $n$ and $T$ control the sensitivity of hand status change, the smaller they are, the more step clips are obtained. In our implementation, we empirically take $n = \frac{FPS}{6}$ and $T = \frac{FPS}{10}$ (FPS: frames per second of the video and $T \leq n$). Figure 6a shows the summation of scores within the window of length $n$ across a whole video as an example. We take down all the moments the hand status switches from 'idle' to 'active' and from 'active' to 'idle' as 'starts' and 'ends' of segments. In Figure 6a, the status change happens at the moment the score summation equals $T$.

Having obtained an initial segmentation $C_l^{initial}$ which is shown on the first subplot of figure 6b, we can see that many segments only last for a relatively short period. This can be caused by the uncertainty when status change and usually happens at the beginning and end of a step. We regard the segments which are less than half a second as a false positive segmentation and remove them directly. The result after short segment filtering $C_l^{filtering}$ is shown on the second subplot of Figure 6b.

### 3.3.2 Clips Re-connection by Similarity Check

Compared with the number of frames per second, the sliding window length and threshold chosen to determine the hand status is relatively small. This makes the hand status prediction more sensitive to ballistic hand movements. At the same time, more incorrect step splits could be introduced. We thus follow a 'break and reconnect' strategy to solve over-segmentation. In other words, we check the similarity of the 'active objects' extracted from two adjacent clips and determined whether to link them up. The clips that have links



(a) Score calculation
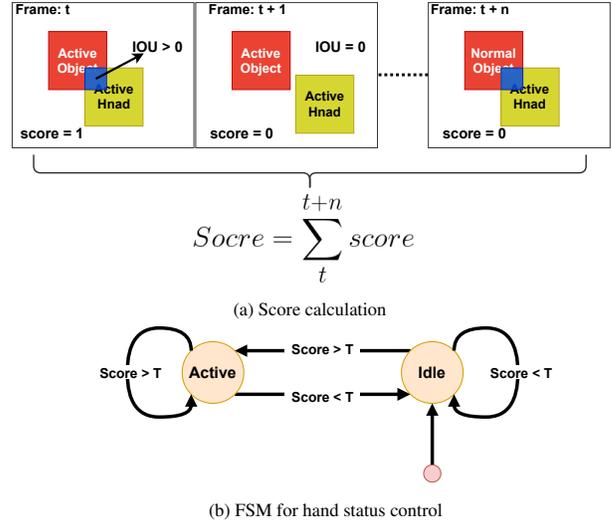


(b) FSM for hand status control

Figure 5: Full hand status determination process

between them are reconnected as one segment.

Recall that the process of determining a hand-object interaction pair in section 3.3.1. If a frame is considered as having an interaction on a given hand, we assume there is only one active object for that hand. For each hand, every frame in $C_l^{filtering}$ has a corresponding object crop. As per 1, the similarity measure averaged image similarities between two adjacent segments by densely calculating the crop-wise similarity. Considering the appearance of the object's crop within the segment is not always the same due to, e.g. pose changes, the similarity measure only happens between the 20% of the ending frames for the foregoing segment and the 20% beginning frames for the posterior segment.

The image-wise similarity is achieved by training a Siamese network [22]. As shown in Figure 7, two images pass through the same feature extractor. For positive samples, we choose an image crop. The training pair is generated by randomly applying image transformation warps and flips. The negative samples come from different object crops. We use *Global Average Pooling* [19] to get the result and *Binary Cross-Entropy Loss* as loss function. Since our HOI dataset has over 20$K$ object instances, it can be used to train our similarity check network.

We use ROC (Receiver operating characteristic) curve analysis to find the best threshold for segment re-connection. We manually select 5 objects and find 10 instances for each of them from different frames in our HOI dataset. For each object class, we run 1000 times similarity check (the positive ground truth rate is set to 50%) by selecting images pairs randomly. With different thresholds from $0 - 1$, we plot the ROC curve shown on the bottom of Figure 7.

### 3.4 Fusion of Two Hand Streams

To get the final task step decomposition segmentation, the results from left and right hands need to be combined into one. We hypothesize that for a single hand-object interaction process, the attention can only be put on one 'active object'-'active hand' pair (there are at most two processes in one frame). We define the attention as 'hand attention', and the hand gained attention as 'principal hand', correspondingly, the hand with less or no attention is called 'secondary hand'.

We hypothesise that the hand attention as defined here is highly related to the relative position of the two hands, especially when both of them are interacting with different objects. To learn the

(a) The 'Score' summation within a sliding window across all frames. $T$ is the threshold of hand status change.



(b) Segmentation results after 'Score' thresholding in (a) (top), the results further optimized by filtering out some segments (middle) and the results after clips re-connection (bottom).
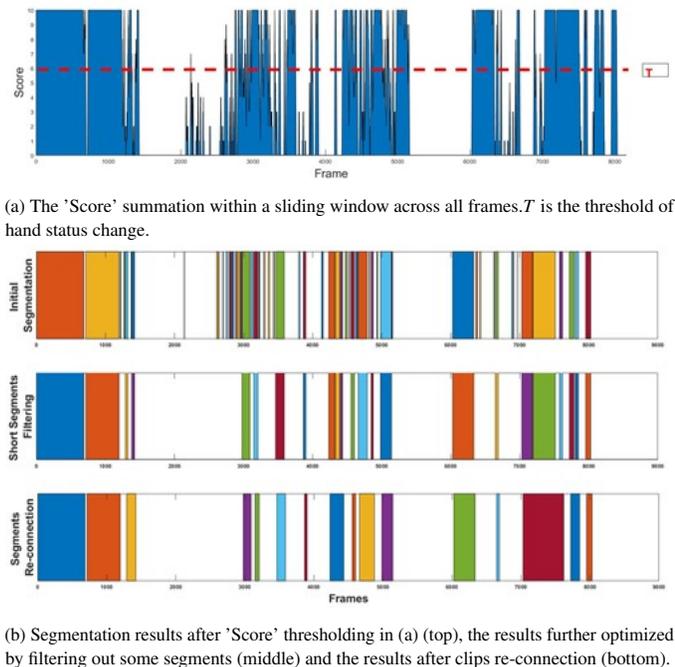
Figure 6: The example of video segmentation process.

distribution characteristics of the hand-based attention mechanism, we follow the rules of determining interaction status to extract all the cases in our FPV-HOI Detection dataset when both hands are in 'interaction with objects'. Then the attention on hand is labelled manually (labelling the hand that dominates the action). We plot the relative position of two hands when they are all recognized as in 'interaction status'. For example, (1) in figure 8 shows the left-hand position distribution when the right hand is normalized to the centre of the hand attention map. From distributions (1) and (2) in figure 8, we observe that the secondary hand in most of the cases is located in a position that is spatially lower to that of the principal hand. While we use multiple and varied hand activities in our training datasets (GTEA [18], Epic-Kitchen [3], FPHA [5] and our Chinese Tea Making datasets). We agree that there is the scope that more data labelled from multiple persons will inform better the distribution of these events.

## 4 EVALUATION

In this section, we use different experiments to (1) evaluate the performance of our hand-object interaction detector quantitatively. (2) quantitatively evaluate the performance of our video segmentation results on GTEA [18]. (3) evaluate our Chinese-tea making based on expert instruction segmentation and participants' annotations.

### 4.1 Dataset Selection

Here we list all the datasets we use for training and evaluation.

- **FPV-HOI**. This is the dataset we collected 3000 images from EPIC-KITCHENS [3], GTEA [18] and FPHA [5] and labelled by us. It contains 6 kinds of labels. We train our HOI detector with this dataset.

- **EPIC-KITCHEN 'unseen'**. EPIC-KITCHEN is an egocentric dataset on kitchen work recording. Action labels are provided, but there are no scripted tasks included which prevent its use in our study here. However we use the images in its 'unseen' kitchen as the test set for our HOI detector.
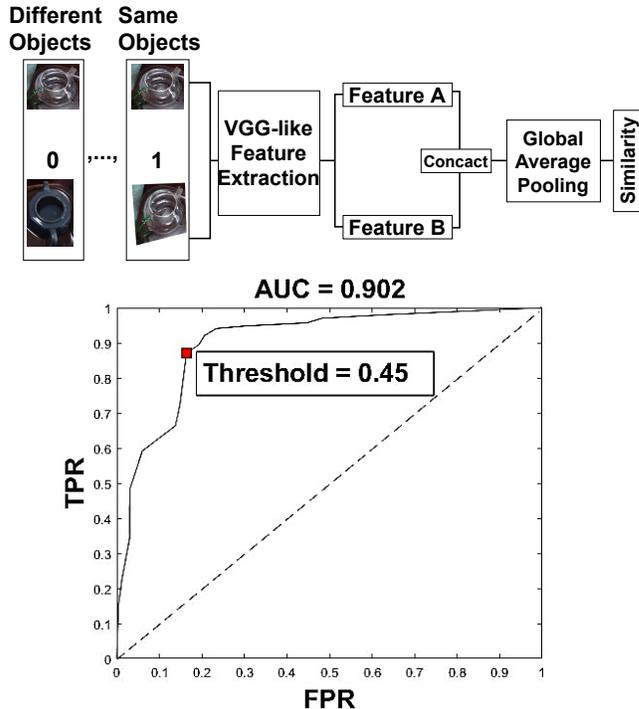


Figure 7: The Siamese net based image wise similarity measure network (top). And the threshold selection for similarity. (bottom)

- **GTEA**. GTEA is a dataset that has 28 egocentric videos of sequential-step task performing. The action annotations are given in a verb-noun form. We use it for our video segmentation evaluation.

- **Chinese Tea Making**. Chinese Tea Making Dataset is an egocentric dataset collected by us. Three experts are invited to perform Chinese tea making twice without any prior knowledge of computer vision. Each video lasts for about 2.5 minutes containing 12 different steps. What makes this dataset complex is the cross hand-object manipulation for each step. Multiple objects and hands may have interaction.

### 4.2 Evaluation of Hand-Object Interaction Detection

There are several hand-object interaction datasets of interest for our work. Such as EPIC-KITCHENS [3], VLOG [4], GTEA [18], 100DOH [27] and AVA [9]. The EPIC-KITCHEN provides auto-labelled object bounding boxes only. In the case of 100DOH, it provides labelling of full-hand-state and interaction objects, but most of the data is third-person based and cannot be distinguished from the available labels. Although our hand-object interaction detector has qualitatively good performance on different datasets (Fig 9), building a general HOI detector is not the scope of this work. To make sure there are enough hand-object interactions in our testing set, we select 300 images from the 'unseen' set (the set not in our training set) of EPIC-KITCHENS and only label with hand status and active objects (compared with training images, no 'normal object' is labelled in a testing set). The results are shown in table 1.

On the testing data, we achieve 97.32% precision (true positive over true positive plus false positive) on hand-object interaction. We use a relatively high threshold for the detector because the accuracy loss caused by false negatives for a single frame can be compensated by using video temporal information via hand status control. The
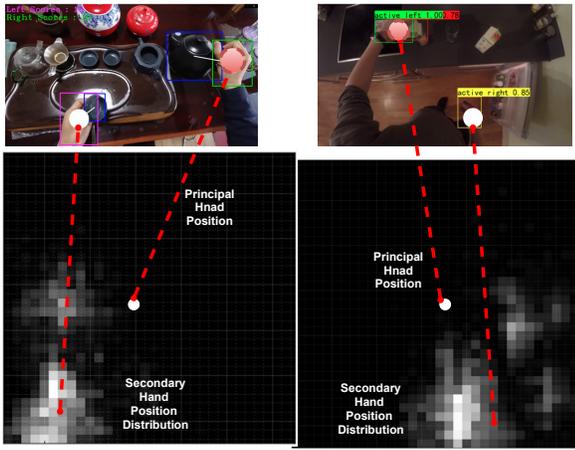
(1)      (2)

(3)

Figure 8: (1) Shows the left-hand position distribution when the right hand is considered as principal hand (the hand dominates the action). (2) shows the right-hand position distribution when the left hand is considered as the principal hand. We take the centre of the hand bounding box as hand position. The position of the right/left hand on (1)/(2) is normalized to the centre of the hand attention map (white dot), and all distances are normalized within $32 \times 32$ (the map size). (3) shows the cases we consider in segments fusion: A, B and C are the cases the IOSA (intersect over the smaller area) of two segments is less than 0.5. D, E and F are the cases IOSA is greater than 0.5, the hand with attention dominates the segmentation.

figure shows example results of our HOI detector. We detect all possible objects and hand-object interactions.

## 4.3   Evaluation on Video Segmentation

We mainly evaluate our task segmentation workflow on GTEA [18] and our Chinese Tea Making dataset since, in contrast to EPIC-KITCHENS, they do have repeated videos for the sequential-step activities. GTEA is a breakfast-making dataset that provides two kinds of labels: action-based and object-based. The action-based label is composed of a verb and a noun. For example, 'take cheese', 'open cheese' and 'put cheese on the bread' are different segments in the action-based label. While in the object-based label, they are all 'cheese' relevant actions and belonging to the same segment. The object-based segmentation has fewer breaks across the video, which is more natural for a human to understand, and we follow this labelling style in GTEA dataset evaluation. Chinese Tea Making dataset is more realistic and thus more complex than GTEA. It is hard to define the ground truth of video segmentation based on either
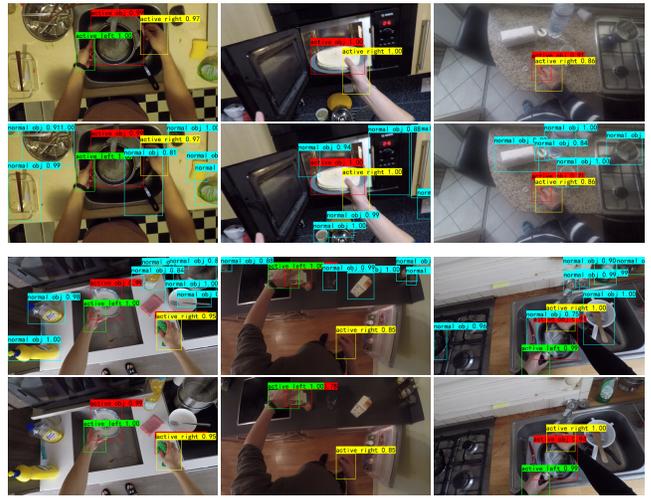


Figure 9: Some examples of hand-object interaction results. Top rows has the results including 'normal object' detection.

Table 1: The results of HOI detector on the 'unseen' subset from EPIC-KITCHENS [3]. In the first row, AH: 'active hand', AO: 'active object' 'HOI': The hand object interaction. The leading number means the number of instances in total. From the second row, TP: number of true positives, FP: number of false positives, TPR: True positive rate. Our HOI detector achieves 97.32% precision which means the possibility a frame is wrongly classified as HOI is low.

| | 385 AH | 322 AO | 385 HOI |
|---|---|---|---|
| TP | 341 | 263 | 291 |
| FP | 45 | 31 | 8 |
| TPR | 88.57% | 81.63% | 75.58% |
| Precision | 88.34% | 89.46% | **97.32%** |

'active objects' or actions because two hands may have different in-hand objects and actions on a timestamp. We use the instruction based annotation (segmented by an expert) and participants' annotation as two distinct references to evaluate our results on the Chinese Tea Making dataset.

We use a segmentation **F1** score with overlapping thresholds (IOU) at 10%, 20% and 50% as quantitative evaluation which is proposed by [15]. As for processing time, our video segmentation currently runs in an offline manner. With an 'i7-6700' CPU laptop and 'Quadro M2000' GPU, the processing time on hand-object detector is about 2 frames per second. We expect this can be optimized in a number of ways.

### 4.3.1   Results on GTEA dataset

As mentioned above, GTEA is a dataset that has action and in-hand object labels with time segments. Because our method is in-hand object-centred, we use the object-based label as ground truth. This level of annotation is also the reason why the comparisons with other works are not performed. We report $F1$ score in table 2.

Figure 11 depicts a typical result of 'cheese sandwich making'

Table 2: The $F1$ score with overlapping thresholds (IOU) 10%, 20% and 50% on GTEA dataset.

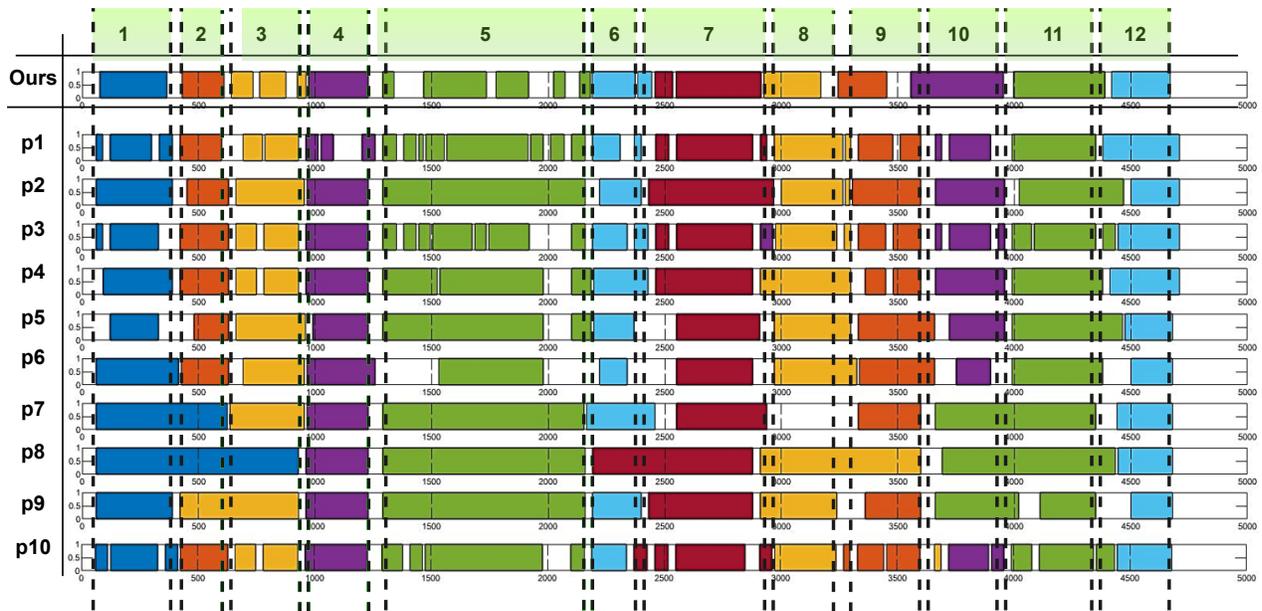| GTEA Results | F1@10% | F1@30% | F1@50% |
|---|---|---|---|
| | 81.67% | 73.01% | 69.50% |

Figure 10: An example segmentation result on a Chinese Tea Making video. It shows the results from our paradigm and the annotations from 10 invited participants. The numbers in top row represent the instruction based expert's segmentation.
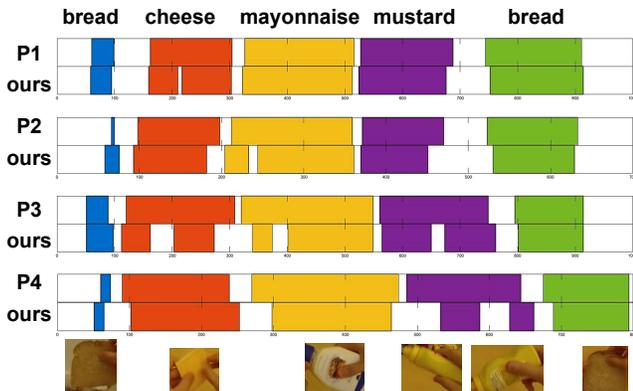


Figure 11: Example result of 'cheese sandwich making' from GTEA.

in GTEA. Our prediction has accurate separation of different in-hand object manipulation and achieves 81.67%, 73.01%, 69.50% $F1$ score, respectively. The corresponding image crops shown in the bottom of 11 come from the object crop of the middle frame in each segment and align the nouns in the ground truth. We trace the reason for having an extra segment in cheese preparation and find it is caused by the failure detection of cheese and the similarity check between the two segments. For the highly deformable object like cheese, the 'cheese' and the 'In-hand cheese' could be two different objects from the perspective of the detector. This could be solved in future work by modelling the difference of an object in different states.

### 4.3.2 Results on Chinese Tea Making Dataset

To have a clearer understanding of this dataset, we list all the objects involved and steps for performing the dataset in table 3.

The steps shown in the table 3 are written and agreed upon by all experts, which summarizes the process of Chinese tea making. However, it is neither action by action nor an object by object instruction.

For example, in the Chinese tea making process, the step 1: 'pour the tea from the kettle to the teapot' is not a step with a single action or a single object. As shown in figure 12, our system roughly split the first step: 'pour the tea from the kettle to the teapot' into 6 stages according to the clips from left and right hand before fusion (the arrow show the clips to be connected). 'Teapot lid', 'teapot', 'kettle' and both hands are involved in this single step. Coincidentally, the left hand's object (teapot lid) does not change through this step. By conducting crop similarity check, all clips from the left hand are reconnected into one. By the rules of segments fusion in section 3.3 and the case $E$ in figure 8. The long segment dominants in video segmentation. However, in another video performed by another expert, the step 1 is implemented by a single right hand, leading to different segmentation results. It is subjective which segmentation approach is better in terms of guidance without further human involved experiments. Thus, the instruction is used as a reference instead of definite ground truth.

In order to have a more objective view of our results, we invite 10 participants to segment into steps the 6 videos in the Chinese Tea Making dataset. They are required to split the videos into different steps. Each step can be expressed with a 'verb' plus 'noun' pair. Gaps are allowed in segmentation. Figure 10 shows our segmentation result and participant's annotations for one video. The numbers from $1 - 12$ represent the instruction from the expert's segmentation, and $p1 - p10$ are the results from the participants. As expected, they show some consensus and some disagreements with expert's segmentation. Like step 3, the actual content is pouring out water from two teacups. It is reasonable to either consider this step as a segment or split it into two. In table 4, we report the $F1$ score of our results on expert's segmentation, average results on participants' annotations and the average results between participants' and expert's. The score shows our result has close parallels with both expert's annotation and participants' annotations.

Our initial motivation for this work was to close the loop and provide video guidance for people to achieve tasks, such as Chinese Tea making. However, due to the COVID-19 restrictions, we cannot conduct experiments with volunteers. Instead, we opted for show a segmentation result with the most segments and invite a new
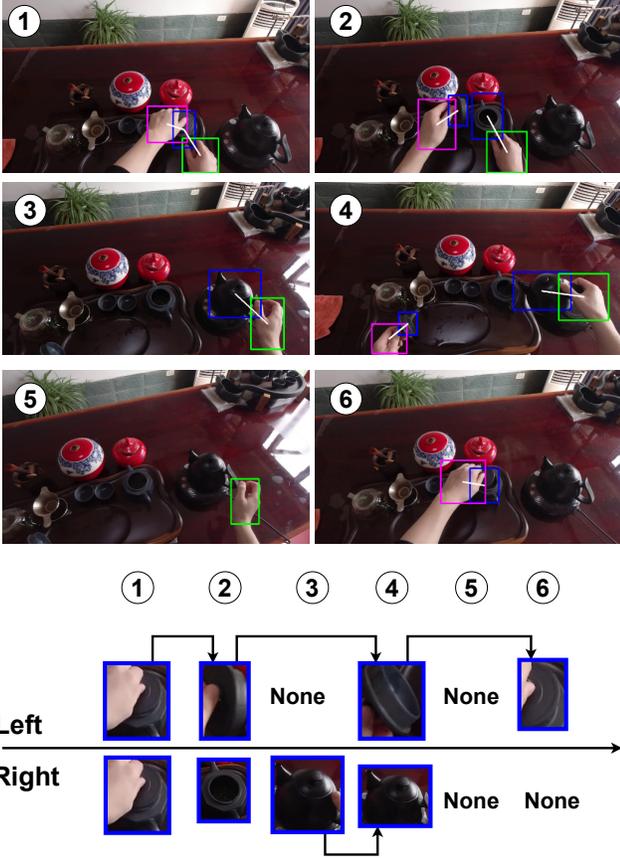
| Chinese Tea Results | F1@10% | F1@30% | F1@50% |
|---|---|---|---|
| ours on instruction | 75.52% | 70.74% | 66.28% |
| ours on 10 labels | 72.34% | 68.30% | 63.95% |
| 10 labels on instruction | 69.19% | 61.76% | 54.55% |

| *Human description based on our results* | |
|---|---|
| 1 | 'Fill the teapot with hot water' |
| 2 | 'pour the hot water into teacup' |
| 3 | 'pour out the water from teacup' |
| 4 | '*not clear*', 'take the brush', '*dip the brush*' |
| 5 | 'take the lid of the jar', '*not clear*', 'take the spoon' |
| | '*not clear*', 'scoop the tea from the jar' |
| | 'remove the lid of the teapot', 'put the tea into teapot' |
| | '*not clear*', 'put the lid back to the jar' |
| 6 | 'pour water to the teapot' |
| 7 | 'put the lid back to the teapot' |
| | 'put the filter onto the glass teapot' |
| | 'filtering the tea' |
| 8 | 'pour the tea into the teacup' |
| 9 | 'pour out the tea from teacup' |
| 10 | 'fill the teapot with hot water' |
| 11 | 'filtering the tea' |
| 12 | 'put back the lid of glass teapot' |
| | 'pour the tea into the teacup' |

Table 5: The description of a participant after watching the segmented video clips. The numbers are the instruction based segmentation which are used as a reference.

participant to transcribe the content of each video clip with verbs and nouns. If they think the clip is not clear enough or cannot express the segment with a verb-noun pair, the clip is labelled as ***not clear***. This approach offers an indication of the quality of the potential guidance and task decomposition.

From the table 5, except for the step 4 and 5, descriptions align well to the instructions. In step 4, the actual step is 'brushing the table' (in order to get rid of the water poured out from the previous step). Due to the false negatives on 'brush' detection, the interaction status is not recognized correctly. While the step 5 is over segmented. Because there are many object transfers between both hands exist in this step, some clips don't show clear intentions.

## 5 CONCLUSION

Task guidance is an essential application for MR systems. In this paper, we focus on the critical competence of task decomposition to support video-based guidance. Our approach uses CNN modules and explainable Finite State Machines to extract the steps that decompose hand activities automatically. We evaluate in real tasks on both public and specifically collected video datasets. Our approach combines hand detection and objects similarity check module to edit videos into steps automatically. We address egocentric step segmentation and segment fusion problems by analysing hand-object interactions and reasoning about hand and object activity within frames. Our results show that we can achieve high precision in step decomposition for unseen tasks, and our results agree within the levels of subjectivity that volunteers can judge. We believe MR systems need to tackle the crucial problem of automated content editing, and methods that develop this aspect for real-life activities will help get us closer to the broader adoption of MR systems. In future work, we will explore methods of delivering the extracted video content to users for task guidance.



Figure 12: The visualization of step 'pouring water form the kettle to the teapot'.

Table 3: Instruction steps in Chinese Tea Making dataset

| Object | Number | Object | Number |
|---|---|---|---|
| teapot | 1 | jar | 7 |
| teapot lid | 2 | jar lid | 8 |
| kettle | 3 | wooden spoon | 9 |
| teacup | 4 | tea filter | 10 |
| brush | 5 | glass teacup | 11 |
| metal tea scoop | 6 | glass teacup lid | 12 |

The steps written with object number and example hand-object interaction case (the hand usage differs in different video):

| Steps | Left Hand | Right Hand |
|---|---|---|
| 1.pouring water form 3 to 1 | 2 | 3 |
| 2.pouring water from 1 into 4 | \ | 1 |
| 3.pouring water out of 4 | \ | 4 |
| 4.brushing the table | \ | 5 |
| 5.scoop the tea from 7 to 1 | 8, 6 | 9, 2 |
| 6.pouring water form 3 to 1 | \ | 3, 2 |
| 7.pouring tea from 1 into 11 | 12 | 10, 1, 11 |
| 8.pouring tea from 11 into 4 | 11 | 11 |
| 9.pouring tea out of 4 | \ | 4 |
| 10.pouring water form 3 to 1 | \ | 3 |
| 11.pouring tea from 1 into 11 | 2, 12, 10 | 1 |
| 12.pouring tea from 11 into 4 | 11 | 11 |

## REFERENCES

[1] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris, and I. Patras. Video summarization using deep neural networks: A survey. *arXiv preprint arXiv:2101.06072*, 2021.

[2] G. Bleser, D. Damen, A. Behera, G. Hendeby, K. Mura, M. Miezal, A. Gee, N. Petersen, G. Maçães, H. Domingues, D. Gorecky, L. Almeida, W. Mayol-Cuevas, A. Calway, A. G. Cohn, D. C. Hogg, and D. Stricker. Cognitive learning, monitoring and assistance of industrial workflows using egocentric sensor networks. *PLOS ONE*, 10(6):1–41, 06 2015.

[3] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018.

[4] D. F. Fouhey, W.-c. Kuo, A. A. Efros, and J. Malik. From lifestyle vlogs to everyday interactions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4991–5000, 2018.

[5] G. Garcia-Hernando, S. Yuan, S. Baek, and T.-K. Kim. First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. *arXiv preprint arXiv:1704.02463*, 2017.

[6] R. Girdhar and K. Grauman. Anticipative video transformer. *arXiv preprint arXiv:2106.02036*, 2021.

[7] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448, 2015.

[8] M. Goto, Y. Uematsu, H. Saito, S. Senda, and A. Iketani. Task support system by displaying instructional video onto ar workspace. In *2010 IEEE International Symposium on Mixed and Augmented Reality*, pp. 83–90. IEEE, 2010.

[9] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6047–6056, 2018.

[10] R. Hanson, W. Falkenström, and M. Miettinen. Augmented reality as a means of conveying picking information in kit preparation for mixed-model assembly. *Computers & Industrial Engineering*, 113:570–575, 2017.

[11] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016.

[12] Y. Huang, Y. Sugano, and Y. Sato. Improving action segmentation via graph-based temporal reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14024–14034, 2020.

[13] S. B. Kang and K. Ikeuchi. Determination of motion breakpoints in a task sequence from human hand motion. In *Proceedings of the 1994 IEEE International Conference on Robotics and Automation*, pp. 551–556. IEEE, 1994.

[14] A. Kis, L. Huber, and A. Wilkinson. Social learning by imitation in a reptile (pogona vitticeps). *Animal cognition*, 18(1):325–331, 2015.

[15] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager. Temporal convolutional networks for action segmentation and detection. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 156–165, 2017.

[16] T. Leelasawassuk, D. Damen, and W. Mayol-Cuevas. Automated capture and delivery of assistive task guidance with an eyewear computer: the GlaciAR system. In *Proceedings of the 8th Augmented Human International Conference*, pp. 1–9, 2017.

[17] S.-J. Li, Y. AbuFarha, Y. Liu, M.-M. Cheng, and J. Gall. Ms-tcn++: Multi-stage temporal convolutional network for action segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[18] Y. Li, Z. Ye, and J. M. Rehg. Delving into egocentric actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 287–295, 2015.

[19] M. Lin, Q. Chen, and S. Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.

[20] Y. Lu and W. Mayol-Cuevas. Higs: Hand interaction guidance system. In *2019 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*, pp. 376–381. IEEE, 2019.

[21] W. W. Mayol and D. W. Murray. Wearable hand activity recognition for event summarization. In *Ninth IEEE International Symposium on Wearable Computers (ISWC'05)*, pp. 122–129. IEEE, 2005.

[22] I. Melekhov, J. Kannala, and E. Rahtu. Siamese network features for image matching. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pp. 378–383. IEEE, 2016.

[23] N. Petersen and D. Stricker. Learning task structure from video examples for workflow tracking and authoring. In *2012 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pp. 237–246. IEEE, 2012.

[24] J. Platonov, H. Heibel, P. Meier, and B. Grollmann. A mobile markerless ar system for maintenance and repair. In *2006 IEEE/ACM International Symposium on Mixed and Augmented Reality*, pp. 105–108. IEEE, 2006.

[25] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779–788, 2016.

[26] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015.

[27] D. Shan, J. Geng, M. Shu, and D. F. Fouhey. Understanding human hands in contact at internet scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9869–9878, 2020.

[28] D. Wang, D. Hu, X. Li, and D. Dou. Temporal relational modeling with self-supervision for action segmentation. *arXiv preprint arXiv:2012.07508*, 2020.

[29] Z. Wang, Z. Gao, L. Wang, Z. Li, and G. Wu. Boundary-aware cascade networks for temporal action segmentation. In *European Conference on Computer Vision*, pp. 34–51. Springer, 2020.

[30] J. Zubizarreta, I. Aguinaga, and A. Amundarain. A framework for augmented reality guidance in industry. *The International Journal of Advanced Manufacturing Technology*, 102(9):4095–4108, 2019.