

---

This is an electronic reprint of the original article.  
This reprint may differ from the original in pagination and typographic detail.

Gökce, Zeynep; Pehlivan, Selen

## Temporal modelling of first-person actions using hand-centric verb and object streams

*Published in:*  
SIGNAL PROCESSING: IMAGE COMMUNICATION

*DOI:*  
[10.1016/j.image.2021.116436](https://doi.org/10.1016/j.image.2021.116436)

Published: 01/11/2021

*Document Version*  
Peer-reviewed accepted author manuscript, also known as Final accepted manuscript or Post-print

*Published under the following license:*  
CC BY-NC-ND

*Please cite the original version:*  
Gökce, Z., & Pehlivan, S. (2021). Temporal modelling of first-person actions using hand-centric verb and object streams. *SIGNAL PROCESSING: IMAGE COMMUNICATION*, 99, Article 116436.  
<https://doi.org/10.1016/j.image.2021.116436>

# Temporal Modelling of First-Person Actions using Hand-Centric Verb and Object Streams

Zeynep Gökce<sup>a,c</sup>, Selen Pehlivan<sup>b</sup>

<sup>a</sup>*Computer Engineering Department, TED University, Ankara, Turkey*

<sup>b</sup>*Department of Computer Science, Aalto University, Espoo, Finland*

<sup>c</sup>*TUBITAK Space Tech. Research Institute, Ankara, Turkey*

---

## Abstract

Analysis of first-person (egocentric) videos involving human actions could help in the solutions of many problems. These videos include a large number of fine-grained action categories with hand-object interactions. In this paper, a compositional verb-noun model including two complementary temporal streams is proposed with various fusion strategies to recognize egocentric actions. The first step is based on construction of verb and object video models as decomposition of actions with a special attention on hands. Particularly, the verb video model that is the spatial-temporal encoding of hand actions and the object video model that is the object scores with hand-object layout are represented as two separate pathways. The second step is the fusion stage to identify action category, where distinct verb and object models are combined to give their action judgments. We propose fusion strategies with recurrent steps collecting verb and object label judgments along a temporal video sequence. We evaluate recognition performances for individual verb and object models; and we present extensive experimental evaluations for action recognition over recurrent-based fusion approaches on the EGTEA Gaze+ dataset.

*Keywords:* first-person vision, egocentric vision, action recognition, temporal models, RNN

---

## 1. Introduction

With the increasing availability and popularity of the wearable cameras, first-person (egocentric) vision offers an interesting scenario to study action recognition problem. Recordings with these cameras have become a part

of daily life and evaluation of actions on these recordings gains decisive importance for applications, in particular, for health monitoring, autonomous driving, robotics and entertainment. For instance, these videos can be analyzed for monitoring the patient activities to detect early signs of dementia [1, 2], or for monitoring the driver’s behavioral status to provide necessary assistance for safe and comfortable driving [3]. In robotic, these kinds of videos are useful to make the robot learn the human motion structure from the first-person view [4]. Besides, tracking and understanding human actions in first-person videos are important for developing feasible virtual reality applications [5].

Unlike third-person setting with fixed camera view, the first-person videos are recorded from the perspective of camera wearer, and usually hands and objects bear the most significant clues to determine action in various scenarios [6]. Due to similar hand-movements and hand-object interactions, these videos include large number of action categories with high inter-class similarities (e.g., *take tomato*, *put tomato*, *mix salad*). In addition to these characteristics, these videos have new challenges such as uneven camera transitions, frequent illumination changes, and limited camera vision. Thus, action recognition task on these videos will be better succeeded with fine-grained evaluation.

Many previous studies focus on modeling egocentric actions as composition of appearance and motion-based features. In this paper, we first decompose the egocentric actions into semantically meaningful and complementary components, verbs and objects [7]. Then, we target the appearance and motion-based features of each component for the purpose of fine-grained analysis and we model their temporal dynamics. Our model aims to process large number of distinct action categories through decomposition and fine-grained analysis while guaranteeing the recognition performance.

Our proposed model is demonstrated in Figure 1 and composes of three main parts, Verb, Object and Fusion models. Particularly, the verb model corresponding to the hand action representation and the object model corresponding to the interactions are represented as two separate pathways which are decomposition of actions. Finally, fusion is the action model employing various fusion strategies based on recurrent neural networks (RNN) on individual verb and object model judgments. Verb model is a verb classifier that takes successive  $C$  clips, each consisting of  $N$  successive frames ( $N=16$  for C3D), and returns the verb scores per clip as an output. Object model is an object detection network taking  $C$  video frames, and returns object

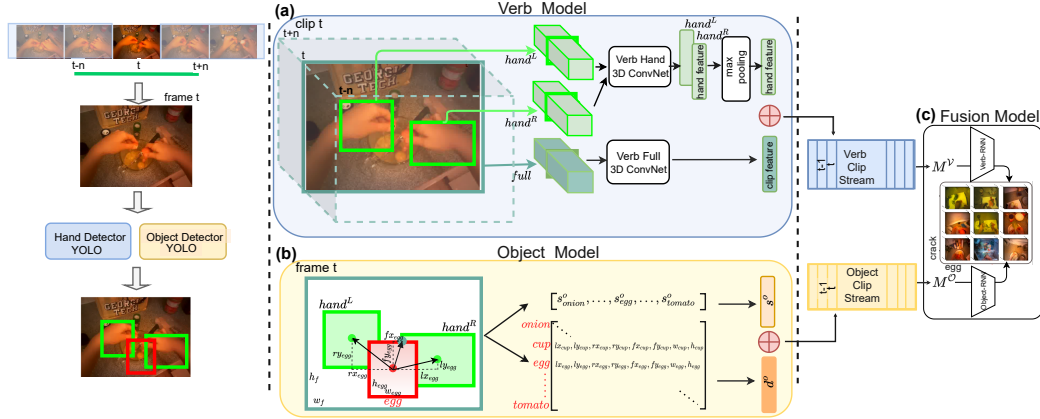


Figure 1: Overview of the proposed first-person action recognition model. Our model proposes a compositional model with two pathways, verb and object streams, respectively, and fusion strategies with recurrent structures. (a) Verb stream possesses the hand-centric action information and (b) the object stream possesses the appearance information of nearby objects with hand-interactions. Then (c) streams are fused using proposed fusion strategies to recognize actions.

proposals.

One of our motivation in the proposed model is to define a hand-centric action decomposition model as hands perform the main action and they are the target of attention in egocentric videos. First, we believe that verbs are strongly related to hand appearance and motion. Thus, our verb model prioritizes the fine-grained analysis of hand actions with spatial-temporal and multiscale modeling in a fully-supervised setting. Model encodes not only hand appearance using RGB-based 3D convolutional neural network (3D ConvNet), but also implicitly hand motion along tubes due to 3D ConvNet (a second stream can be added to encode Flow-based 3D ConvNet). Second, we similarly follow previous studies in mapping nouns to objects, but we believe hand-object interactions is also important to reveal target object. Thus, our noun component encodes two information: object proposals and the spatial layout of hands and object categories in frames.

Other motivation of our work is to contribute to action decomposition with support of recurrent models. Following this, we focus more deeply on modelling the temporal dynamics of hand-actions (verbs) and hand-object interactions (objects) individually and jointly using various fusion strategies. In our verb component, temporal model encodes the hand-centric information across clips. This means that it performs a smoothing over spatial-temporal



multiscale features. In noun component, recurrent model encodes not only the temporal dynamics of object detection scores but also that of spatial layout across frames. This is also a kind of motion-based feature where the model tracks the layout of hands and objects.

Although video-based action decomposition is achieved priority in [7], our design of components is different and specialized over hands. [7] associates verbs with motion-based low-level features without any notion of hands, and represents frames with global dense trajectories [8]. Besides, nouns are associated with appearance-based features and encoded similarly over object proposals. Unlike using all objects, they encode only the objects found near hands. In short, our design is different in how we define the verb and noun components of the decomposition modelling and RNN-based fusion strategies that are proposed to combine complementary components while jointly modelling their temporal dynamics.

We are not the first using hand-centric approach for egocentric models. Hand is the main object and its manipulation is a strong clue for many studies [9, 10, 11, 12, 6]. Although these models are not based on decomposition, compared to these studies our architectural components are designed in a different way, where our model strongly relies on 3D ConvNet architecture for fine-grained evaluation of hand motion, appearance, and its interaction. Unlike attention based models using gaze estimation [13, 14], we examine regions around hand proposals as the target of attention.

The main contribution of our study is that (i) we propose a new first-person action decomposition model with hand-centric verb and noun components. Because hands serve as the target of attention, hand-centric spatial-temporal verb model and interaction feature are proposed. It has been shown that the combination of detectors in multiple scales and interaction feature increase the performance. As another contribution, (ii) we propose various fusion strategies with extension to recurrent structures. We observe that recurrent late fusion strategies outperform early fusion with many architectural advantages. We also show that our decomposition model performs comparable in recognition to conventional action recognition models, but with many architectural advantages. Then, (iii) we populate the recognition model with full supervision, since the action space is significantly large and contains extremely similar action categories. For that purpose, we train the background models of our system on a recent dataset, called EGTEA Gaze+ [13] with additional annotations. In addition to the available version, we gather on videos (1) the actual frames with the action label, and on these frames (2)

the bounding boxes of hands with verb labels (annotations consist of additional verb labels such as *hold*), and (3) the bounding boxes of target objects. Finally, (iv) we present extensive experimental evaluation over fusion approaches on the EGTEA Gaze+ dataset with comparable results.

## 2. Related Work

Within the development of wearable cameras, a wide range of studies are proposed for action recognition in first-person videos [15, 16, 17, 18, 19, 20, 21, 13]. We explore traditional and deep learning based models including appearance-based, motion-based and hybrid models.

### 2.1. Appearance Based Models

Appearance cues related to objects, hands and gaze are informative for first-person videos and they are used in many studies.

Objects have crucial knowledge to describe first-person actions by revealing human-object interactions. In many studies [16, 9], hand-object interactions are modelled over region of interests (ROIs) to understand egocentric activities. According to Fathi et al.[9], the target object is generally visible in the center of video frames and a new model with two steps is proposed over the GTEA dataset. The first step is segmenting videos into foreground and background regions using optical flow, SIFT, color histograms. Foreground segments are further decomposed into hands and active objects. The second step performs recognition using Multiple Instance Learning (MIL) over object segments. According to Fathi et al. [16], fine-grained actions are categorized using hand interaction features (such as optical flow of both hand and object, hand pose, hand location, hand size, and left/right-hand relative location). As their previous work [9], hands, foreground objects and background are segmented and Adaboost [22] classifier is used for recognition. The proposed model results in an accuracy of 45% on the GTEA dataset. Another object-based model by Pirsiavash and Ramanan [17] is a temporal pyramid based model provided to define the usage of the objects in videos for action recognition. HOG [23] features are used for object modelling. With linear SVM classifier, performance is achieved up to 77% using object information over the ADL dataset. Recently, Cartas et al. [24] present another object-based model with two steps over the GTEA dataset. First, the hand region is segmented to get object region in frames using Multiscale Combinational Grouping method [25]. Then, a star-structured region model, R\*CNN [26],

136 is used for more than one region classification. Last, the output of R\*CNN  
137 as contextual cue is given to LSTM to predict action category.

138 Gaze information is another cue for first-person action recognition, since  
139 camera wearer generally focuses on the point where the action is performed.  
140 Visual features extracted around the gaze regions are more informative com-  
141 pared to features extracted on other regions. Fathi et al. [27] extend object-  
142 based model of their previous work [16] with addition of gaze appearance.  
143 The SVM classifier is used for action categorization using object-based, gaze-  
144 based appearance features. Object-based features are extracted from object  
145 classifiers including object context around the gaze point. Using gaze infor-  
146 mation that is given with the GTEA Gaze+ dataset, the unrelated back-  
147 ground objects are eliminated and 47% performance is achieved compared to  
148 27% of [16] on the same dataset without using gaze information. Similarly,  
149 Li et al. [28] develop a model for gaze prediction in first-person videos using  
150 hand/head movement, hand location and hand pose.

## 151 2.2. Motion Based Models

152 First-person videos capture motion information from camera wearer’s  
153 head, hand and eye movements. Besides object-based models, which is known  
154 as appearance-based models, motion-based models are also proposed to rec-  
155 ognize the first-person actions in the literature. Kitani et al. [29] model  
156 motion in first-person sport activities using motion histograms. The motion  
157 histograms are based on optical flow of the scene. Due to the unsupervised  
158 scenario, Dirichlet process mixture models are proposed to get action cate-  
159 gories using the motion histograms. Li et al. [30] model motion information  
160 using Dense Trajectories [31] as a baseline descriptor.

## 161 2.3. Hybrid Models

162 Appearance and motion domains are composed (i.e., stream-based models  
163 [32]), since fusing them is more informative for first-person action recogni-  
164 tion [6, 33, 34, 35].

165 Ma et al. [6] model object appearance and motion information as a two-  
166 stream network. The first stream analyzes appearance in three steps; seg-  
167 mentation, localization and object recognition successively, and the second  
168 stream analyzes motion using optical flow features. Yansong Tang et al. [33]  
169 propose a tri-stream network that integrates depth knowledge besides ap-  
170 pearance and motion information and test over RGB-D egocentric dataset  
171 (THU-READ). Action prediction is calculated by taking the average score

of three streams. Hahn et al. [34] propose a model using visual information from videos and textual information from recipe of these videos as well. The proposed model has three steps which are action proposal, object recognition, and recipe alignment steps. In action proposal step, video frames are localized with Bidirectional LSTM with two classes, action or not-action. In object recognition step, ResNet101 network [36] is trained for object classification along frames having actions according to the action proposal step. Finally, in the recipe alignment step, the action category is predicted using NLP model.

Recently, G. Kapidis et al. [35] introduce a multi-modal approach based on sequential learning to recognize egocentric actions on EPIC-Kitchens dataset [37]. LSTM is trained over feature sequence which consists of frame-based hand coordinates as motion knowledge and presence of object as appearance knowledge. Hand and objects are detected by YOLOv3 [38]. In contrast, our action model is trained over clip-based verb scores based on hand regions as motion information and detected object scores as appearance information using recurrent models.

### 3. Proposed Approach

Many studies emphasize the importance of appearance and motion-based lower-level modelling in action video understanding either with simple concatenation [34, 35] or with stream-based structures [32, 6, 33]. In this study, we aim at stream-based structures, but in semantic-level as decomposing actions into two complementary pathways that are the verb and the object components. Our model targets temporal modelling on each pathway with a strong attention to hands that perform actions to manipulate surrounding objects. The verb component is defined as a hand-centric temporal model. The model consists of short-term temporal modelling of hand regions within each video clip in multiple-scales using spatial-temporal 3D ConvNet models [39, 40]; and long-term temporal modelling of hands over clips of videos using RNN model. Complementary to this, the object component is defined as based on a temporal modelling of objects and hand interactions. The model consists of objects extracted by YOLOv2 object detector [41] and interaction-based spatial layout features, and further relies on long-term modelling of temporal dynamics of hand-object interactions over frames of videos using RNN model. Our aim is to perform action recognition by combining pretrained verb and object models using various fusion strategies.

### 208 3.1. Temporal Modelling of Verbs

209 Assuming the camera wearer’s attention is on the hands in first-person  
210 videos, the verb model is hand-centric that it models short term temporal dy-  
211 namics of hands performing action (hands perform *take* during *take tomato*)  
212 in multiple scale. Our verb model is fully-supervised and trained using clips  
213 including verb-labelled mid-frames and hands on these frames. It composes  
214 of two sub-models in different scales, namely full-scale and hand-scale verb  
215 models. With clips including hands, while the full-scale verb model repre-  
216 sents coarse-grained verb description of the clips by utilizing whole frame  
217 region (covering hands, objects and scene); the hand-scale verb model repre-  
218 sents fine-grained verb description by utilizing zoomed regions around hands.  
219 Given a video, each sub-model extracts stream of clip features per video. Par-  
220 ticularly, these streams correspond to  $V \times C$ -dimensional verb matrices with  
221  $V$ -dimensional features and  $C$  clips (When softmax outputs are used as clip  
222 representation,  $V$  corresponds to verb categories. Otherwise,  $V$  corresponds  
223 to the dimension of intermediate layer). These matrices are further combined  
224 over scales into a single verb matrix as the video verb representation. Figure  
225 1 (a) shows an overview of the verb stream.

#### 226 3.1.1. Full-Scale Verb Representation

227 The purpose of the full-scale verb model is to encode coarse-level verb  
228 category details of video clips. Model is based on 3D ConvNet architecture  
229 C3D [39] which takes the video clips as inputs and produces the category  
230 score vectors as outputs. Given a video with  $C$  successive clips extracted  
231 using temporal stride of two (i.e., dropping every other clip), each clip is  
232 embedded into a  $V$ -dimensional feature by the full-scale verb model. As the  
233 video full-scale hand representation, the model returns  $V \times C$ -dimensional  
234 verb matrix.

235 Following the original setting of [39], the 16-frame video clips are ex-  
236 tracted and each is resized to  $112 \times 112 \times 3 \times 16$  before fed into C3D <sup>1</sup>. Unlike  
237 the original C3D model [39] trained over randomly selected video clips, our  
238 model is trained over ground truth video clips that include verb-labelled mid-  
239 frame with hand performing verb (see Section 4.1 for the dataset details).

---

<sup>1</sup><https://github.com/hx173149/C3D-tensorflow>

### 240 3.1.2. Hand-Scale Verb Representation

241 As part of a hand-centric approach, the purpose of the hand-scale verb  
 242 model is to encode fine-level verb category details of video clips. Focusing on  
 243 hand actions, this model utilizes hand regions instead of looking at a video  
 244 in full-scale mode. It consists of two parts: hand detector and verb classifier.  
 245 Given a video with clips, first a hand detector localizes hand regions in mid-  
 246 frames of these clips, and then a spatial-temporal verb classifier takes the  
 247 hand-volumes around these regions to identify verb categories.

248 To localize hands in video frames, we train a hand detector using the state-  
 249 of-the art object detector YOLOv2 [41]. The original YOLO architecture is  
 250 fine-tuned for binary classification of hands (hand/not-hand) on our hand  
 251 dataset gathered from the EGTEA Gaze+ dataset (see Section 4.1 for the  
 252 dataset details). During training of the hand-scale verb model, the YOLO  
 253 hand detector is used to obtain hand-volumes, where volumes are particularly  
 254 tubes cropped around hand proposals containing the hand-region in its mid-  
 255 frame and lasts 16 frames. Hand proposals having 0.5 overlap with verb-  
 256 labelled hands on the ground truth frames are used to train our hand-scale  
 257 verb model.

258 Following the full-scale model, the hand-scale one is also based on 3D  
 259 ConvNet architecture C3D [39] which takes hand-volumes as inputs and pro-  
 260 duces category scores as outputs. Given a video with the same set of  $C$  clips,  
 261 hand-volumes are computed from hand proposals on mid-frames of these clips  
 262 and fed into 3D ConvNet model to get volume features. As the video hand-  
 263 scale verb representation, the model returns  $V \times C$ -dimensional verb matrix.  
 264 Here, the cropped hand-volume is also resized into  $112 \times 112 \times 3 \times 16$  before  
 265 fed into the 3D ConvNet. Detection of multiple hand-volumes on the same  
 266 video frame is possible, since some action categories are performed by both  
 267 hands (e.g., while *open* and *take* are performed by one hand, *cut* and *mix* are  
 268 performed by two hands). In that case, the second hand acts as an auxiliary  
 269 hand to help the main hand performing action. For example, in verb *cut*,  
 270 one hand cuts the object while the other holds the object. Following this,  
 271 hand-scaled verb model is trained with one extra verb category, verb *hold*,  
 272 that is also a ground truth label of hand regions in our dataset. In case of  
 273 having multiple hands on a frame, such as *cut*, we apply max-pooling over  
 274 features of hands to reduce into a single feature vector.

### 275 3.1.3. Video Verb Recognition

276 Given a video, our model extracts multiple verb matrices in multiple  
277 scales, full-scale and hand-scale, as verb representations of the video. These  
278 matrices are combined into a single verb matrix  $M^V$  using max-pooling. Verb  
279 matrix is particularly a sequence of clip features. Having training videos with  
280 various number of clips, we introduce a count-based verb model (see Section  
281 4.3.1) using histograms and RNN-based verb model (see Section 4.4.1) to  
282 recognize the verb category of videos.

### 283 3.1.4. Discussion on 3D ConvNet Architectures

284 In this study we use two variants of 3D ConvNet architectures, C3D  
285 network [39] and I3D RGB network [40] to encode verbs in spatial-temporal  
286 domain.

287 C3D network has 5 convolutional and 5 pooling layers, where each con-  
288 volution layer is immediately followed by a pooling layer. Then, network  
289 includes 2 fully connected layers and a linear classifier to predict action cate-  
290 gories [39]. Compared to C3D, I3D is a denser network of inception-v1 with  
291 3 convolution layers, 9 inception modules and  $7 \times 7$  average-pooling layer pre-  
292 ceding the last linear classifier [40]. Two-stream 3D ConvNet extension is  
293 available with I3D RGB and I3D Flow pathways, but we use RGB modality  
294 to encode our verbs.

295 C3D network is shallower than I3D network. But the fact that it is shal-  
296 lower makes the network appealing to train for each scale. Following this, we  
297 first concentrate on C3D network and train two models for hand-scale and  
298 full-scale keeping the original setting. These models are used to extract low  
299 level clip and hand features (see Section 3.1.1 and Section 3.1.2). Later, I3D  
300 RGB network is used in final experiments to evaluate the performance im-  
301 provement. Instead of training two I3D models for each scale, we fine-tune the  
302 model with a simple network over our dataset. Selected intermediate layers  
303 of pre-trained I3D RGB network are used to extract coarse-level features for  
304 full-scale model and fine-level features for hand-scale model. Later, a simple  
305 shallow network is trained for classification over concatenated full-scale and  
306 hand-scale verb representations (please see Section 4.4.4 for network details).

### 307 3.2. Temporal Modelling of Objects

308 Objects manipulated by hands is used to model noun component of our  
309 decomposition model. First, we aim to find which objects appear in the  
310 video, and then we encode object information to reveal their interactions

311 with hands. Figure 1 (b) indicates an overview of the object stream. Our  
 312 object model is trained over object bounding boxes that are annotated on  
 313 verb-labelled frames with hand annotations (see Section 4.1 for the dataset  
 314 details). Given a video, object model extracts stream of frame features per  
 315 video. Particularly, these streams correspond to  $O' \times C$ -dimensional verb ma-  
 316 trices with  $O'$ -dimensional features and  $C$  clips.

### 317 3.2.1. Object Representation

318 To localize objects in video frames, the object detector YOLOv2 [41] is  
 319 fine-tuned with object categories gathered from the EGTEA Gaze+ dataset  
 320 (see Section 4.1 for details). During training, the YOLO object detector is  
 321 used to obtain object proposals on video frames extracted using temporal  
 322 stride of two frames. Later, each frame is encoded with a  $O'$ -dimensional  
 323  $[\mathbf{s}^o, \mathbf{d}^o]$  feature vector, where  $\mathbf{s}^o$  is an  $O$ -dimensional object score vector,  
 324 and  $\mathbf{d}^o$  is a  $8 \times O$ -dimensional object distance vector computed to represent  
 325 hand-object interactions.

326 Proposals having 0.4 overlap with objects on the ground truth verb-  
 327 labelled frames are extracted and used for training our object model. Having  
 328 detected proposals, max-pooling is applied to pool over the confidence scores  
 329 of the detected objects of the same category, and this results in a score vec-  
 330 tor  $\mathbf{s}^o = [s_1^o, \dots, s_i^o, \dots, s_O^o]$ , where  $s_i^o$  is the maximum score over all detected  
 331 proposals belonging to category  $i$  and  $O$  is the number of object categories.

332 In addition to score vector, interactions of detected objects with detected  
 333 hands (see Section 3.1.2 for detected hands) are encoded using a distance-  
 334 based representation  $\mathbf{d}^o = [\mathbf{d}_1^o, \dots, \mathbf{d}_i^o, \dots, \mathbf{d}_O^o]$  to encode spatial layout, where  
 335  $O$  is the number of object categories and  $\mathbf{d}_i^o$  is the distance vector belonging  
 336 to category  $i$  as follows,

$$\begin{aligned}
 \mathbf{d}_i^o &= [lx_i, ly_i, rx_i, ry_i, fx_i, fy_i, w_i, h_i] \\
 lx_i &= (cx_i - cx_{lhand})/w_f, \quad ly_i = (cy_i - cy_{lhand})/h_f \\
 rx_i &= (cx_i - cx_{rhand})/w_f, \quad ry_i = (cy_i - cy_{rhand})/h_f \\
 fx_i &= (cx_i - cx_f)/w_f, \quad fy_i = (cy_i - cy_f)/h_f \\
 w_i &= w/w_f, \quad h_i = h/h_f
 \end{aligned} \tag{1}$$

337 where  $lx_i$  and  $ly_i$  show the scaled x-distance and y-distance between the  
 338 center of the left-hand and the center of the object  $i$ , respectively.  $rx_i$  and  
 339  $ry_i$  represent the scaled distances between the center of the right-hand and  
 340 the center of the object  $i$ .  $fx_i$  and  $fy_i$  show x-distance and y-distance of



object center  $i$  to frame center. Variables  $w_f$ ,  $h_f$ ,  $w$  and  $h$  are for the frame-width, frame-height, width and height of detected object  $i$ , respectively.

Detected hands are categorized as a left-hand or a right-hand based on their relative distances. The detected hand whose center is closer to the top left corner of the frame is classified as the left-hand and the right-hand otherwise. If there is only one hand detected, we duplicate the values for both hands. If multiple proposals of the same object category are detected, the proposal having a minimum Euclidean distance with any hand is selected (If there is no proposal for an object category, we insert zero values.).

As the video representation, the model returns a  $(O+8\times O)\times C$ -dimensional object matrix  $M^O$  as stacked  $[\mathbf{s}^o, \mathbf{d}^o]$  features over video frames.

### 3.2.2. Video Object Recognition

Given a video, our model extracts a matrix  $M^O$  as object representation of the video. Object matrix is particularly a sequence of frame features. Having training videos with various number of frames, we introduce a count-based object model (see Section 4.3.3) using histograms and RNN-based object model (see Section 4.4.2) to recognize the object category of videos.

### 3.3. Temporal Modelling of Actions as Fusion of Verb-Object Pairs

Fusion is the last step of our proposed model to combine verb and object streams for recognizing actions. Since action videos consist of sequence of short-term clips, modeling of temporal relations between consecutive clips are important for action recognition. Such temporal modelling is also critical to smooth information over clips. Given the decomposition of actions into verb and object streams per video, we introduce multiple strategies with early and late fusion techniques using recurrent neural network (RNN) models for encoding temporal dynamics of actions. In this section, we first introduce a new count-based baseline model, and then we describe five different fusion strategies over verb-object streams for recognizing actions. Figure 2 shows the proposed fusion strategies. The architectural details of the best performed neural network for the proposed fusion strategies are given in Table 4 and experimental evaluations are reported in Table 7 and Table 8 .

**Count-based verb-object multiplication.** This model has no recurrent step and no training for action recognition (see Figure 2 (a)). Verb and object category scores for videos are obtained using the convolutional neural network

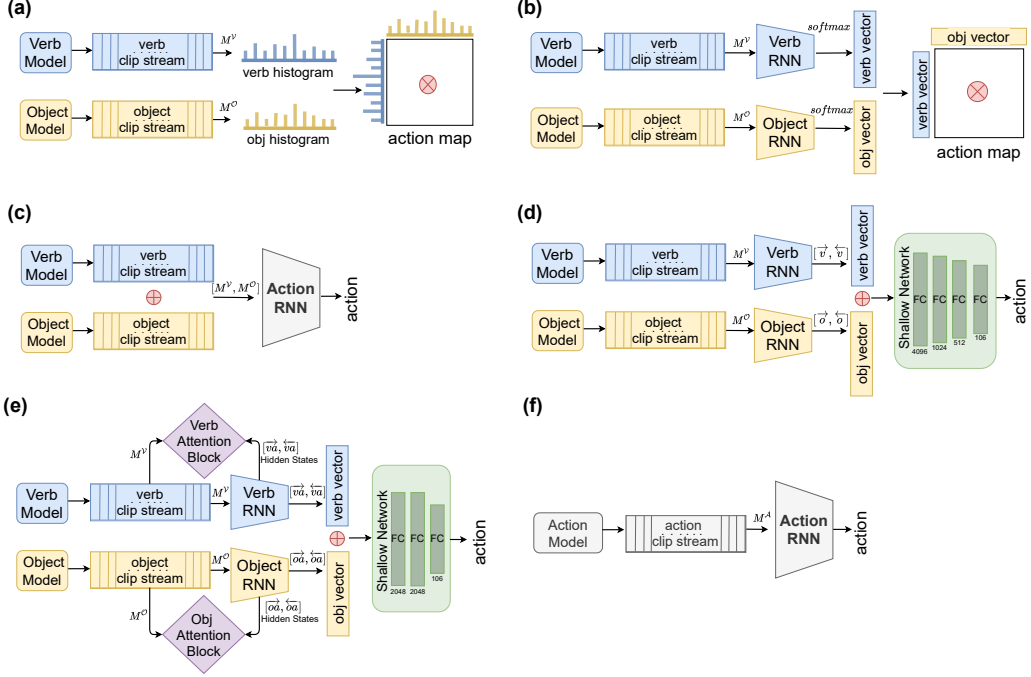


Figure 2: Fusion strategies for action recognition. (a) Count-Based Verb-Object Multiplication Baseline model, (b) Recurrent Verb-Object Multiplication model, (c) Recurrent Verb-Object Early Fusion model, (d) Recurrent Verb-Object Late Fusion model, (e) Recurrent Verb-Object Attention Late Fusion model, (f) Recurrent Action Baseline model (verb vector is either “scores” or “features”, since softmax scores and intermediate level features are used interchangeably in the experiments).

softmax predictions per clip, and the model performs a simple multiplication for recognition.

To compute verb stream, count-based multiplication model uses 3D ConvNet softmax prediction scores over clips of video samples. Given a video, two  $V \times C$ -dimensional verb matrices with  $V$  verb categories and  $C$  clips are extracted using full-scale and hand-scale 3D ConvNet verb models, and these matrices are combined into a single score matrix by max-pooling (see Section 3.1). Computing verb stream, each clip is assigned to a verb category with the maximum score over  $V$  categories. Later, a histogram showing the distribution of verb categories over  $C$  video clips is computed and L1 normalization is applied to eliminate the effect of video length. To sum up, the video is represented as a  $V$ -dimensional verb category score vector  $\mathbf{v}$ . Similarly, to compute object stream, objects are detected in frames using the object

389 detection model (this fusion strategy does not use spatial layout feature,  $\mathbf{d}^o$ ,  
 390 see Section 3.2.1). Then, each frame is assigned to an object category with  
 391 the maximum detection score. Later, a histogram showing the distribution  
 392 of object categories over  $C$  video frames is extracted and L1 normalization  
 393 is applied to eliminate the effect of video length. To sum up, the video is  
 394 represented as a  $O$ -dimensional object category score vector  $\mathbf{o}$ .

395 Inspiring from a recent study on human-object interactions in still im-  
 396 ages [42], we combine verb vector  $\mathbf{v}$  and object vector  $\mathbf{o}$  using a simple  
 397 multiplication as follows,

$$\begin{aligned} A &= \mathbf{v} \cdot \mathbf{o}^T \\ A' &= A \odot B \end{aligned} \tag{2}$$

398 where  $A$  is a  $V \times O$ -dimensional estimation map with scores for all action  
 399 categories corresponding to combinations of all *verb-object* category pairs.  $B$   
 400 is a  $V \times O$ -dimensional ground truth binary mask where 1 shows the existence  
 401 of a *verb-object* category pair, 0 shows the nonexistence of the pair in the  
 402 dataset (e.g., *cut-fridge* pair is 0 since it is not an action in our dataset).  
 403 In order to evaluate the scores of the subset of *verb-object* pairs existing in  
 404 the dataset, the estimation map  $A$  is masked by binary mask  $B$  and the final  
 405 result is matrix  $A'$ .

406 Finally, the *verb-object* pair with the maximum value over matrix  $A'$  is  
 407 assigned as the predicted category of the given test sample. Particularly, this  
 408 fusion strategy returns the prediction of action categories without training.

409  
 410 **Recurrent verb-object multiplication.** This model uses the score vec-  
 411 tors of RNN-based verb and object models as inputs (see Section 3.1.3 and  
 412 Section 3.2.2 ), and action recognition stage is performed with a simple mul-  
 413 tiplication over these score vectors without training (see Figure 2 (b)). First,  
 414 two individual RNN models are trained over verb matrix  $M^v$  and object  
 415 matrix  $M^o$  of training videos, namely verb-RNN and object-RNN models.  
 416 During testing, the verb-RNN returns the  $V$ -dimensional verb category score  
 417 vector  $\mathbf{v}$  ( $V$  is the number of verb categories) and the object-RNN returns  
 418  $O$ -dimensional object category score vector  $\mathbf{o}$  per video sample ( $O$  is the num-  
 419 ber of verb categories). Then, we similarly apply multiplication and masking  
 420 using Eq. 2. Please note that this fusion strategy does not use spatial layout  
 421 feature,  $\mathbf{d}^o$ , and object-RNN is trained over stream of score vectors  $\mathbf{s}^o$  (see  
 422 Section 3.2)

423

424 **Recurrent verb-object early fusion.** This model includes a single recur-  
 425 rent model to recognize the action category of the video (see Figure 2 (c)).  
 426 Extracting verb (see Section 3.1) and object (see Section 3.2.1) matrices,  
 427  $M^V$  and  $M^O$ , we concatenate them as a video representation. Then, a single  
 428 RNN model, action-RNN, is trained to predict action categories.

429

430 **Recurrent verb-object late fusion.** Similar to second fusion strategy, Re-  
 431 current Verb-Object Multiplication, two individual recurrent neural network  
 432 models, verb-RNN and object-RNN, are trained over verb and object streams  
 433 on training video samples (see Figure 2 (d)). Then, the score vectors from the  
 434 RNN models,  $\mathbf{v}$  and  $\mathbf{o}$ , are concatenated as a video representation. As an ad-  
 435 ditional training stage, a shallow network with a set of fully connected layers  
 436 is trained to predict action categories over concatenated representation.

437 We also extend this model with forward and backward feature vectors  
 438 when BiLSTM recurrent models are available. Let  $[\vec{\mathbf{v}}, \overleftarrow{\mathbf{v}}]$  be the concatena-  
 439 tion of forward and backward-direction BiLSTM recurrent function outputs  
 440 of the verb network, and  $[\vec{\mathbf{o}}, \overleftarrow{\mathbf{o}}]$  be forward and backward-direction BiLSTM  
 441 recurrent function outputs of the object network. Then, these vectors from  
 442 the verb and object recurrent models are concatenated and  $[\vec{\mathbf{v}}, \overleftarrow{\mathbf{v}}, \vec{\mathbf{o}}, \overleftarrow{\mathbf{o}}]$  is  
 443 used as a video representation. Later, a shallow network with a set of fully  
 444 connected layers is trained to predict action categories.

445

446 **Recurrent verb-object attention late fusion.** This model includes an  
 447 additional attention module to encode temporal information over RNN mod-  
 448 els (see Figure 2 (e)). On each component of our two-stream model, a re-  
 449 current neural network layer is trained with a self-attention module [43] as  
 450 follows,

$$\alpha^V = \text{softmax}(\mathbf{w}_2 \tanh(W_1 M^V)), W_1 \in \mathbb{R}^{256 \times V}, \mathbf{w}_2 \in \mathbb{R}^{1 \times 256} \quad (3)$$

451 where  $M^V$  is the verb stream matrix. The attention block takes  $M^V$  and  
 452 BiLSTM outputs as input. The  $M^V$  is fed into a fully connected layer with a  
 453  $256 \times V$ -dimensional weight matrix  $W_1$ , followed by  $\tanh()$  function. Then, a  
 454 score vector is produced by applying a vector of parameters  $\mathbf{w}_2$  that is  $1 \times 256$ .  
 455 Later, we sum up the RNN model hidden states according to the weight  
 456 provided by  $\alpha^V$  to get a vector representation  $\mathbf{a}^V$ . Here, vector  $\alpha^V$  represents

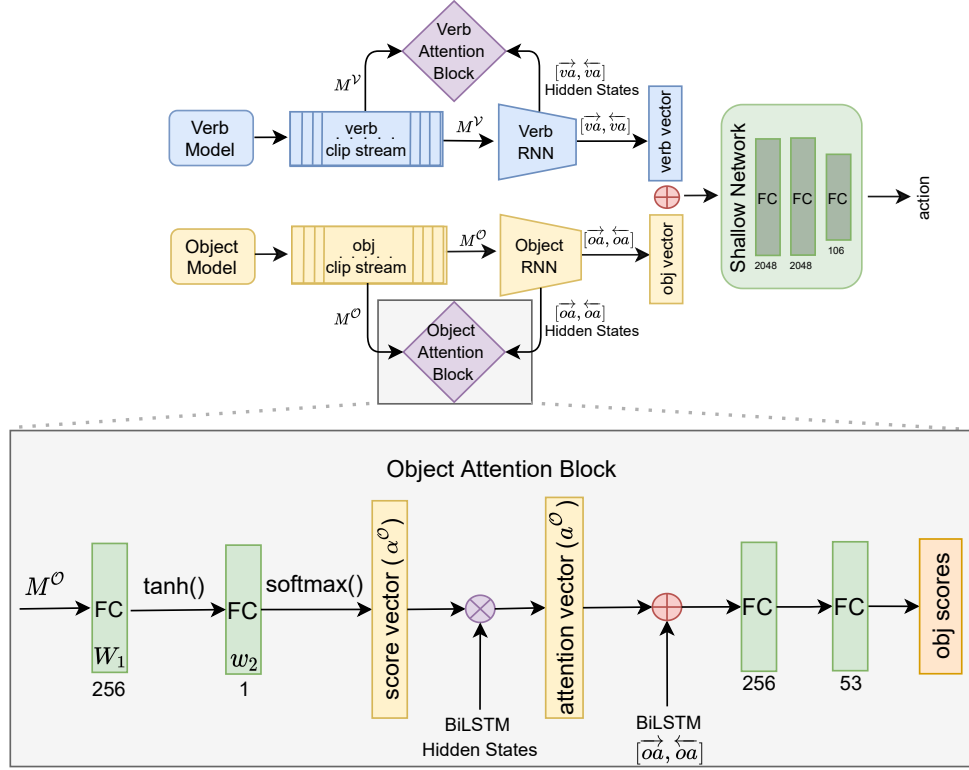


Figure 3: Recurrent Verb-Object Attention Late Fusion strategy. This model extends late fusion strategy with an additional attention block that can be applied both to verb and object streams. Attention block helps to encode temporal context within an attention vector.

the temporal attention of the video. The softmax() function ensures all the computed weights sum up to 1.

The object stream is similarly trained over  $M^O$  using RNN with attention module. Assuming BiLSTM as our RNN structure, forward and backward function outputs of verb-BiLSTM trained with attention block,  $[\vec{v}a, \overleftarrow{v}a]$ , and forward and backward function outputs of object-BiLSTM,  $[\vec{o}a, \overleftarrow{o}a]$ , are concatenated as a video representation. Later, a shallow network with a set of fully connected layers is trained to predict action categories. For the verb and object streams, we have added attention block as seen in Figure 3.

**Recurrent action baseline.** This model is a baseline model using 3D ConvNet architecture trained over action categories (see Figure 2 (f)). Given

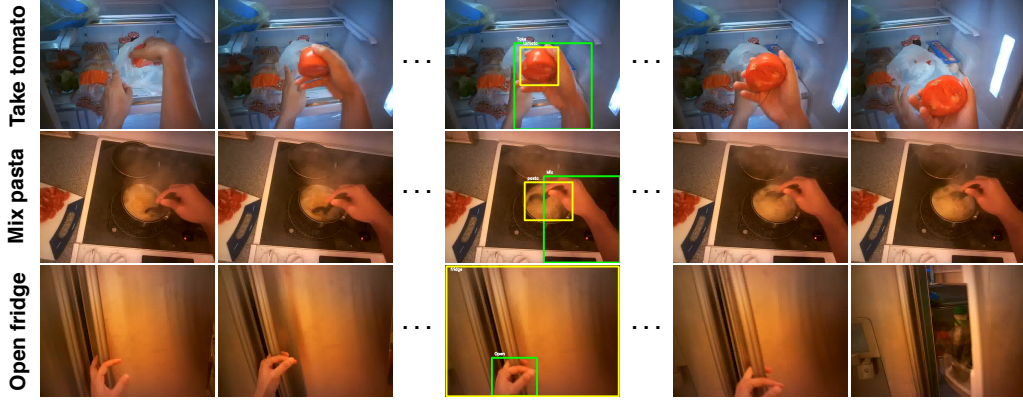


Figure 4: Our sample annotations on the EGTEA Gaze+ dataset. We annotate video frames, hands and the objects having consistent motion with the action in training videos.

469 a video with  $C$  clips, 3D ConvNet architecture C3D takes clips as inputs and  
 470 returns category scores as outputs. Video is represented as a matrix of clip  
 471 features,  $M^A$ . Then, an RNN model, action-RNN, is trained over  $M^A$  to  
 472 predict action categories.

## 473 4. Experimental Evaluation

### 474 4.1. EGTEA Gaze+ Dataset with Frame Level Annotations

475 We perform our experiments on the EGTEA Gaze+ dataset [13, 44]  
 476 which includes first-person meal preparation activity videos. The dataset  
 477 is extended from the GTEA Gaze+ and it consists of 86 cooking videos per-  
 478 formed by 32 different subjects. It includes 106 fine-grained action categories  
 479 with 19 verb and 53 object categories. While some action categories (e.g.,  
 480 *cut tomato*) refer one object category, others (e.g., *pour water-faucet-pot*)  
 481 refer multiple objects. The dataset contains 3 train-test splits, each of which  
 482 has 8229 training and 2022 test video samples. Since the study [13] pub-  
 483 lished the EGTEA Gaze+ dataset reports results on split1, we perform the  
 484 experimental evaluation of models on the same split.

485 Current version of the EGTEA Gaze+ dataset does not include frame  
 486 level annotations. In our study, video frames, the hands and the objects  
 487 having consistent motion with the action category are annotated in training  
 488 videos (see Figure 4). The annotated frames with hands correspond to the

middle frames of clips used in training of 3D ConvNet verb models (see Section 3.1). *hand* is not included among object categories, therefore *wash hand* action category in *wash* verb and *hand* object is just recognized using additional *wash-hand* verb category (no object). Moreover, some video samples cannot be assigned to a verb category due to the absence of hands, therefore *inspect-read recipe* action is identified using only *recipe* object. As a result, annotations for 20 verb categories are provided for training of verb models (19 verbs + *hold*).

Moreover, objects interacting with hands are used in our object model. Therefore, the frames with hand annotations are labelled with object locations and same object categories are populated with additional annotations from other frames of the videos. Annotations of *water*, and *seasoning* are skipped since they are quite ambiguous to label. Moreover, the *hand* is not included among object categories, since the hand annotations are used to train auxiliary hand model in the hand-scale verb model (see Section 3.1.2). As a result, annotations for 50 out of 53 object categories are provided for the object models.

#### 4.2. Experimental Setting

The details of the best performed neural network action recognition architectures with their verb and object submodels are summarized in Table 4. For each component of the decomposition model, we have a base model and an RNN model. Base models for verb component are build based on C3D and I3D networks, and base model of object component is based on YOLO. While we keep the original architectures of C3D and YOLO, we fine-tune the I3D architecture with a simple shallow network over our dataset. Later, our Verb and Object-RNNs are trained and their outputs are used as inputs of action models. We train the networks for minimizing the cross-entropy loss and we use ADAM as the optimization algorithm.

Late fusion strategy is proposed as our main action recognition architecture with higher recognition performance. The best late fusion model using C3D- $d^o$  (*Verb-Object Late Fusion-fc7(with  $d^o$ )-BiLSTM Bw+Fw* in Table 7) with 50.29% mAcc consists of a shallow network of 3 fully connected layers with 4096, 1024, 512 neurons and a 106-dimensional softmax output (see row 4 in Table 4). It is trained using batch size 10 and learning rate 1e-5. The Verb-RNN model of C3D- $d^o$  (*C3D full+hand (concatenation)-fc7-BiLSTM*, see Table 5) is constructed using 1 layer BiLSTM with 768 cell size and

525 trained using batch size 10, learning rate 1e-5. On the other hand Object-  
 526 RNN model of C3D- $d^o$  (*YOLO-s^o + d^o-BiLSTM*, see Table 6) is constructed  
 527 with 1 layer BiLSTM with 1024 cell size, and it is trained using batch size  
 528 10 and learning rate 1e-4.

529 Later, our model is improved with I3D feature (*Our I3D-full+hand-*  
 530 *BiLSTM Bw+Fw* in Table 8). The I3D action model (see row 5 in Table  
 531 4) consists of 2 fully connected layers with 1024 neurons. It is trained with  
 532 batch size 20 and learning rate 1e-4. The Verb-RNN of I3D model is basically  
 533 constructed using 1 layer BiLSTM with 512 cell size and trained using batch  
 534 size 10 and learning rate 1e-6. The Object-RNN of I3D is same with C3D- $d^o$ .

535 Finally, we extend our late fusion strategy with an attention module, I3D-  
 536 Att as our best performing model (*Our I3D-full+hand-BiLSTM+Attention*).  
 537 The verb and object-RNNs of the model are constructed using 1 layer with  
 538 512 and 256 cell sizes, respectively (see row 6 in Table 4). Training is per-  
 539 formed with batch size 10 and learning rate 1e-4 for verb-RNN and 1e-5 for  
 540 object-RNN, respectively. The action model using Recurrent Verb-Object  
 541 Attention Late Fusion setting contains two 2048-dimensional fully connected  
 542 layers with dropout rate of 0.5. It is trained using batch size 20, learning  
 543 rate of 1e-5.

544 One main problem we face is the imbalanced samples of the dataset.  
 545 Using all ground truth verb data, we observe that the model tends to clas-  
 546 sify categories with more samples, but misclassify the categories with fewer  
 547 samples (e.g., *take* verb category contains 1886 samples, but *squeeze* verb cat-  
 548 egory has 29 samples). Thus, as the accuracy gets increased, the mean class  
 549 accuracy drops significantly. To solve this problem, we train verb models  
 550 using balanced subsets. If a category has more than 250 samples, we ran-  
 551 domly select 250 samples. If a category has less than 200 samples, we apply  
 552 some aggregation strategies and we generate new clips by sampling around  
 553 the ground truth clips. For all other categories, we keep their original num-  
 554 ber of samples. Thus, we train all reported verb recognition experiments  
 555 with balanced data having maximum 250 and minimum 200 samples for all  
 556 categories.

### 557 4.3. Experiments on Count-Based Baseline Models

558 In this section, recognition performances of verb, object and action models  
 559 are evaluated, respectively, when a simple count-based strategy is followed  
 560 to compute video category score over clip categories for the related task.  
 561 For individual verb and object recognition tasks, verb and object streams



Verb Base Model	mAcc (%)
C3D full	39.99
C3D hand20	37.13
C3D hand10	35.91

Table 1: Experiments to show the effect of hand-volume ROI scales on video verb recognition results with mean class accuracies (mAcc).

are the stacked softmax prediction scores over clips and the final video verb category or object category is identified using a simple distribution. For action recognition task, we first obtain the distribution of verb and object categories over clips and we apply a fusion model called Count-Based Verb-Object Multiplication strategy (see Figure 2 (a)).

#### 4.3.1. Count-Based Verb Models in Multiple Scale

Region of interests (ROIs) in different scales encode different amount of information from the hand and the background. In order to examine the effect of scale in hand-centric verb model performance, we evaluate verb recognition in various scales. Verb recognition is analyzed by computing the main verb category appearing in the video. The category with the maximum value on the  $V$ -dimensional verb category score vector  $\mathbf{v}$  (histogram) returns the predicted verb class of the video (see Count-Based Verb-Object Multiplication model in Section 3.3).

The performance of each scale is shown in Table 1. The first verb model is trained in full-scale mode (see Section 3.1.1). Other models, hand-scale verb models, are trained using different enlargement scales around hand bounding boxes (hand10 and hand20 verb models mean that 10% and 20% of enlargements with respect to the width and the height of the detected hand regions are applied, respectively). We observe that the full-scale verb model outperforms the hand-scale verb models, since the full region encodes the scene information including hand motion, object and background. The context provided by each element enhances the ability to recognize action in videos [16].

#### 4.3.2. Combination of Count-Based Verb Models

When the verb recognition accuracies are investigated in the category level, it has been observed that the hand-scale verb model outperforms the full-scale one in some verb categories such as *open*, *put*, *crack* verbs. Therefore, we also analyze the performance of verb models when multiple scales

Combination	Verb Base Model	mAcc (%)
<b>weighted average</b>	C3D full+hand10	46.99
	C3D full+hand20	45.63
	C3D full+hand10+hand20	46.38
<b>max-pooling</b>	C3D full+hand10	46.91
	C3D full+hand20	44.91
	C3D full+hand10+hand20	43.19

Table 2: Experiments to show the effect of combined verb models on video verb recognition results with mean class accuracies (mAcc). The set of weight parameters for C3D full+hand10 combination is  $\{\beta_{full} = 0.5, \beta_{hand10} = 0.5\}$ , for C3D full+hand20 combination is  $\{\beta_{full} = 0.5, \beta_{hand20} = 0.5\}$ , and for C3D full+hand10+hand20 combination is  $\{\beta_{full} = 0.4, \beta_{hand10} = 0.3, \beta_{hand20} = 0.3\}$ .

are combined. We use two methods to combine, weighted average and max-pooling. In the first method, softmax values of the clip sequences from full-scale and hand-scale models are weighted averaged at the clip level. The weight parameters  $\{\beta_{full}, \beta_{hand10}, \beta_{hand20}\}$  are empirically searched in the range of  $[0-1]$ . In the second method, the max-pooling is applied over the softmax values of the clip sequences.

Combining verb video representations, we compute the verb category score vector  $\mathbf{v}$ . The category with the maximum value on vector  $\mathbf{v}$  returns the predicted verb class of the video (see Count-Based Verb-Object Multiplication model in Section 3.3). The combination performances are reported in Table 2. The result shows that the combination of softmax values enhances the mean class accuracy (mAcc) of the verb model up to 46.99% (since the full-scale verb model achieves the best accuracy, we keep the full-scale verb model and combine it with the hand-scale models). We observe that the combination of verb models enables the model to capture low-level appearance features both in coarse scale and fine scale. Particularly, hand-scale and full-scale models vote for the clip score together where hand-scale model returns scores over local hand ROI and this helps to recognize harder video instances that are misclassified without local details. The weights of the best combination in weighted averaging method are 0.5 and this shows the best score is achieved with equal contribution of the scales. Although the accuracy of the weighted average method is slightly higher than the max-pooling method, we select the max-pooling method for the count-based action recognition experiments due to its simplicity.

Action Model with Simple Counting	mAcc (%)
Count-Based Verb-Object Multiplication	33.87
Count-Based Action Baseline	23.89

Table 3: Action recognition results with mean class accuracies (mAcc) using Count-Based Verb-Object Multiplication model and Count-Based Action Baseline model on split1.

#### 616 4.3.3. Count-Based Object Model

617 We also perform the evaluation of object stream at the video level. For  
618 the object recognition, the category with the maximum value on the O-  
619 dimensional object category score vector  $\mathbf{o}$  (histogram) returns the predicted  
620 object class of the video (see Count-Based Verb-Object Multiplication model  
621 in Section 3.3). Particularly, this is to compute how often the main object  
622 category appears in the video and the main object is the assigned category  
623 for a video frame with the maximum YOLO confidence. 63.41% mean class  
624 accuracy (mAcc) is achieved for video object classification on split1 (see row  
625 2 of Table 6). For *trash*, *mixture*, *condiment* object categories, low accuracies  
626 are achieved since these objects are hard to detect. Moreover, we observe that  
627 for some categories labelling is ambiguous and ground truth labeling causes  
628 low accuracies for these object categories. For instance, *tomato container*  
629 instances are visually similar to *grocery bag* instances and they are getting  
630 mixes up with each other. In another example, *fridge*, *fridge drawer* and  
631 *drawer* instances are getting confused and used interchangeably in labelling.

632 We also perform a simpler evaluation of object recognition over video  
633 samples as a baseline, where we do not use the histograms to identify the  
634 predicted object category but instead we find the category with the maximum  
635 score over all YOLO detections of all frames. This results in 50.70% mAcc  
636 value as reported in the first row of Table 6. It shows that histogram based  
637 model better evaluates the object recognition over videos.

#### 638 4.3.4. Count-Based Action Model

639 Our first fusion strategy for action recognition is based on a simple mul-  
640 tiplication of the verb category score vector  $\mathbf{v}$  and the object category score  
641 vector  $\mathbf{o}$  (see Count-Based Verb-Object Multiplication model in Section 3.3).  
642 The verb scores are computed using C3D network. The results are reported  
643 in Table 3 with a mean class accuracy (mAcc) of 33.87% on split1.

644 For comparison, we also construct a Count-Based Action Baseline model  
645 having the same implementation with the verb model, but architecture learns

646 action categories rather than verb categories. The Action Baseline model is  
 647 trained over video clips in a supervised setting using annotated action frames  
 648 (clips). For each clip of the test video, a softmax output over action labels  
 649 is retrieved from the C3D model, then the frequently observed action label  
 650 on the clip sequence of the video is evaluated using a histogram proposed  
 651 in Section 3.3 (instead of applying count based model on verb and object  
 652 streams separately, we apply the same model only on action scores). We ob-  
 653 serve that the Count-Based Verb-Object Multiplication model outperforms  
 654 the Count-Based Action Baseline model with almost 10% accuracy. Although  
 655 multiplication of stream scores is a simple technique for action recognition,  
 656 its performance is higher than the baseline model without any learning (see  
 657 Table 3). Here, we have the same training instances for all models, but verb-  
 658 object multiplication trains over a smaller number of categories with more  
 659 samples for verb and object streams compared to action baseline model. This  
 660 might help neural network models to better train. Moreover, our model con-  
 661 tains fine-grained representations of the video instances compared to action  
 662 baseline model using detection models used in the background for objects  
 663 and hands.

#### 664 4.4. *Experiments on Recurrent-Based Models*

665 In this section, recognition performances of verb, object and action mod-  
 666 els are evaluated, respectively, when recurrent structures are used for the  
 667 related task. Particularly, we evaluate the recurrent-based fusion strategies  
 668 we propose for action recognition (see Figure 2 (b-e)). Recurrent models get  
 669 the verb and the object representations of video as a set of clip features,  
 670 and then they use RNN models (LSTM or BiLSTM) to encode the temporal  
 671 dynamics within streams.

##### 672 4.4.1. *Recurrent-Based Verb Models*

673 We examine the individual performances of recurrent-based verb models.  
 674 From the verb experiments of count-based verb model (see Section 4.3.2), we  
 675 know that the combination of features in multiple scales improves the perfor-  
 676 mance of verb recognition. Therefore, we combine verb models in different  
 677 scales: the full-scale and hand-scale verb models. We use hand20 verb model  
 678 for split1 due to its performance reported in Table 1. Here, verb matrices  
 679 of multiple verb models are combined by concatenation or max-pooling, and  
 680 then the combined feature matrix is fed into either BiLSTM or LSTM re-  
 681 current model to analyze the verb recognition performance. Two different

Fusion Model		Verb Base	Verb-RNN	Object Base	Object-RNN	Action Model
<b>Multiplication C3D soft-max/fc7</b>	soft	C3D Full&Hand original	BiLSTM 100 cell, 1 layer	YOLO original	BiLSTM 100 cell, 1 layer	
	fc7	✓	768 cell, 1 layer	✓	✓	
	input	Clip(112×112×3×16)	full <i>softmax</i> /fc7 hand <i>softmax</i> /fc7	Img(416×416×3)	Obj $s^o$	
	output	Verb-full <i>softmax</i> Verb-hand <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	
<b>Early C3D(concat) softmax</b>	soft	C3D Full&Hand original		YOLO original	BiLSTM 100 cell, 1 layer	
	input	Clip(112×112×3×16)		Img(416×416×3)	Verb-concat <i>softmax</i> Obj <i>softmax</i>	
	output	Verb-full <i>softmax</i> Verb-hand <i>softmax</i>		Obj <i>softmax</i>	<i>softmax</i> (106)	
<b>Late C3D softmax/fc7</b>	soft	C3D Full&Hand original	BiLSTM 100 cell, 1 layer	YOLO original	BiLSTM 100 cell, 1 layer	Shallow Network fc(512)-drop(0.5)- fc(256)-drop(0.5)
	fc7	✓	768 cell, 1 layer	✓	✓	✓
	input	Clip(112×112×3×16)	full <i>softmax</i> /fc7 hand <i>softmax</i> /fc7	Img(416×416×3)	Obj $s^o$	Verb-RNN <i>softmax</i> Obj-RNN <i>softmax</i>
	output	Verb-full <i>softmax</i> Verb-hand <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	<i>softmax</i> (106)
<b>Late C3D bw+fw softmax/fc7 / fc7 with <math>d^o</math></b>	soft	C3D Full&Hand original	BiLSTM 100 cell, 1 layer	YOLO original	BiLSTM 100 cell, 1 layer	Shallow Network fc(1024)-drop(0.5)- fc(1024)-drop(0.5)- fc(512)-drop(0.5)- fc(256)-drop(0.2)
	fc7	✓	768 cell, 1 layer	✓	✓	fc(1024)-drop(0.5)- fc(512)-drop(0.5)
	w/ $d^o$	✓	✓	✓	1024 cell, 1 layer	fc(4096)-drop(0.5)- fc(1024)-drop(0.5)- fc(512)-drop(0.5)
	input	Clip(112×112×3×16)	full <i>softmax</i> /fc7 hand <i>softmax</i> /fc7	Img(416×416×3)	Obj $s^o/s^o + d^0$	Verb-RNN $Bw+Fw$ Obj-RNN $Bw+Fw$
	output	Verb-full <i>softmax</i> Verb-hand <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	<i>softmax</i> (106)
<b>Late I3D bw+fw with <math>d^o</math></b>	model	Our I3D(fine-tuned) fc1(1024)-drop(0.5)- fc2(1024)-drop(0.5)	BiLSTM 512 cell, 1 layer	YOLO original	BiLSTM 1024 cell, 1 layer	Shallow Network fc(1024)-drop(0.5)- fc(1024)-drop(0.5)
	input	<i>Mixed-5c</i> <i>MaxPool3d-5a-2x2</i>	Our I3D <i>fc2</i>	Img(416×416×3)	Obj $s^o + d^0$	Verb-RNN $Bw+Fw$ Obj-RNN $Bw+Fw$
	output	Verb <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	<i>softmax</i> (106)
<b>Late I3D-Att bw+fw with <math>d^o</math></b>	model	Our I3D(fine-tuned) fc1(1024)-drop(0.5)- fc2(1024)-drop(0.5)	BiLSTM-Att 512 cell, 1 layer Att Module	YOLO original	BiLSTM-Att 256 cell, 1 layer Att Module	Shallow Network fc(2048)-drop(0.5)- fc(2048)-drop(0.5)
	input	<i>Mixed-5c</i> <i>MaxPool3d-5a-2x2</i>	Our I3D <i>fc2</i>	Img(416×416×3)	Obj $s^o + d^0$	Verb-RNN $Bw+Fw$ Obj-RNN $Bw+Fw$
	output	Verb <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	<i>softmax</i> (106)
	output	Verb <i>softmax</i>	Verb-RNN <i>softmax</i>	Obj <i>softmax</i>	Obj-RNN <i>softmax</i>	<i>softmax</i> (106)

Table 4: Architectures details of action models from Table 7 and Table 8 with their verb and object base models and RNN structures. ✓ marks indicate the repetition of the above model structure.

Verb Base Model	Feature	Verb-RNN	mAcc (%)
C3D full	softmax	BiLSTM	46.49
C3D full+hand (concatenation)	softmax	BiLSTM	53.95
C3D full+hand (max-pooling)	softmax	BiLSTM	50.43
C3D full+hand (max-pooling)	softmax	LSTM	48.13
C3D full	fc7	BiLSTM	53.23
C3D full+hand (concatenation)	fc7	BiLSTM	<b>59.50</b>

Table 5: Video verb recognition results with mean class accuracies (mAcc) on split1 with recurrent-based methodologies.

intermediate features of C3D verb models are experimented on to construct verb feature matrices, softmax prediction scores and fc7 layer features.

We first test using softmax prediction scores. The verb matrices are the stacked softmax outputs per clip. According to the experimental results given in Table 5, it is observed that the combination of full-scale and hand-scale verb models helps in verb recognition with 7.46% improvement over full-scale verb model on split1. The feature concatenation method outperforms the max-pooling in verb models. Moreover, the verb model with BiLSTM structure (50.43%) gives higher accuracy than LSTM structure (48.13%), and we therefore continue with BiLSTM for the rest of the experiments. When compared to Table 2, it is clearly seen that recurrent verb models outperform the simple count-based models in recognition. This shows that recurrent models are better to model verb streams. The best BiLSTM verb model is constructed using 1 layer with 100 cell size (see Table 4).

In order to improve verb recognition performance in recurrent-based models, we also test our experiments with fc7 layers instead of softmax scores. Extracted fc7 layers from the full-scale and the hand-scale verb models (please note that hand features are combined into a single feature vector using max-pooling, if there are multiple hands) are concatenated per video clip. Then, the video representation as stacked clip features is fed into the BiLSTM model. The experiments are conducted over split1, and we obtain 53.23% accuracy using full-scale verb model and 59.50% accuracy using concatenated features. Results show that fc7 that is an earlier feature layer improves recognition rates significantly. The best BiLSTM verb model with fc7 features is constructed using 1 layer with 768 cell size (see Table 4).

Object Base Model	Feature	Object Recognition	mAcc (%)
YOLO	$s^o$	Max-pooling	50.70
YOLO	$s^o$	Count-based	63.41

Object Base Model	Feature	Object-RNN	mAcc (%)
YOLO	$s^o$	BiLSTM	70.59
YOLO	$s^o+d^o$	BiLSTM	<b>70.83</b>
YOLO	$s^o$	LSTM	68.73

Table 6: Video object recognition results with mean class accuracies (mAcc) using count-based and recurrent-based strategies. For the recurrent models,  $d^o$  indicates that distance-based spatial layout features are also integrated while modelling object features. Otherwise,  $s^o$  features are used. (Please see Section 3.2.1)

#### 707 4.4.2. Recurrent-Based Object Models

708 We evaluate individual performances of recurrent-based object models  
709 using BiLSTM and LSTM structures (see Section 3.2.1). Object features,  
710  $[s^o, d^o]$ , are extracted over video frames by the object model in a matrix  
711 form, and then they are fed into RNN object model. BiLSTM object model  
712 using stand-alone  $s^o$  achieves 70.59% accuracy on split1. The BiLSTM model  
713 has 1 layer with 100 cell size (see Table 4). Extending model with spatial  
714 layout feature using  $[s^o, d^o]$ , BiLSTM object model achieves 70.83% accuracy  
715 over split1. The model has 1 layer with 1024 cell size (see Table 4).

716 According to Table 6, it is observed that the recurrent-based object mod-  
717 els outperform the count-based object model with more than 7% improve-  
718 ment. This means that modeling temporal dynamics for video object recog-  
719 nition significantly improves performance. It has been seen that BiLSTM  
720 object model also improves the accuracy compared to LSTM. However, the  
721 effect of  $d^o$  is very low, but we will later show its effect on action results.

#### 722 4.4.3. Recurrent-Based Action Models

723 We make use of the outputs of the recurrent-based verb and object models  
724 and we train action models based on different fusion strategies (see Table 7).

725 Comparing various fusion strategies, we show simple multiplication re-  
726 sults in comparable performance with more complicated fusion strategies  
727 without any training. Moreover, we show late fusion strategies can achieve  
728 better recognition compared to early fusion with many training advantages.  
729 Moreover, the results of the late fusion strategies outperform other strategies  
730 with addition of bw+fw features and spatial layout feature  $d^o$ . Finally, com-  
731 paring softmax and fc7 results, we observe that RNN achieved comparable  
732 performance on simple softmax features.

Fusion Model	Verb Base Model	Feature	Verb/Object/Action-RNN	mAcc (%)
Verb-Object Multiplication	C3D full+hand	softmax	BiLSTM	45.29
Verb-Object Early Fusion	C3D full+hand	softmax	LSTM	44.36
Verb-Object Early Fusion	C3D full+hand	softmax (max-pooling)	LSTM	44.10
Verb-Object Early Fusion	C3D full+hand	softmax	BiLSTM	45.18
Verb-Object Early Fusion	C3D full+hand	softmax (max-pooling)	BiLSTM	46.47
Verb-Object Late Fusion	C3D full+hand	softmax	BiLSTM	45.45
Verb-Object Late Fusion	C3D full+hand	softmax	BiLSTM Bw+Fw	49.01
Verb-Object Action Baseline	C3D full	softmax	BiLSTM	25.36
Verb-Object Multiplication	C3D full+hand	fc7	BiLSTM	45.54
Verb-Object Late Fusion	C3D full+hand	fc7	BiLSTM	46.62
Verb-Object Late Fusion	C3D full+hand	fc7	BiLSTM Bw+Fw	48.46
Verb-Object Late Fusion	C3D full+hand	fc7 (with $d^o$ )	BiLSTM Bw+Fw	<b>50.29</b>

Table 7: Action recognition results with mean class accuracies (mAcc) on videos using recurrent-based fusion strategies. Bw means Backward function output and Fw means Forward function output. Concatenation is applied to merge full and hand scale verb features unless specified as ax-pooling.  $d^o$  indicates that distance-based features are also integrated while modelling object features. Otherwise,  $s^o$  features are used.

**Recurrent verb-object multiplication.** In this fusion setting, action category of a video is simply identified with multiplication of verb-RNN and object-RNN score vectors (see Figure 2 (b)). The verb category score vector  $\mathbf{v}$  is extracted using the verb model *C3D full+hand(concatenation)-softmax-BiLSTM* (see row 2 in Table 5 and row 1 in Table 4) and object vector  $\mathbf{o}$  is extracted using the object model *YOLO- $s^o$ -BiLSTM* (see Table 6 and row 1 in Table 4), respectively. Given a test video, we simply multiply the verb and the object vectors. The verb-object pair with the maximum value of the matrix obtained by multiplication is selected as the predicted action category of the video. This experiment is applied over split1 with a 45.29% mAcc. We observe that temporal action model significantly outperforms simple Count-Based Verb-Object Multiplication model with 11.42% gain in accuracy (see Section 4.3.4 and Table 7). This means that modelling the temporal dynamics helps in action recognition as well.

Similar experiment is conducted using fc7 features for the verb model *C3D full+hand(concatenation)-fc7-BiLSTM* (see row 6 in Table 5 and row 1 in Table 4) with the same object model. We obtain 45.54% accuracy with a slight improvement over softmax score features.

**Recurrent verb-object early fusion.** In this fusion setting, action recognition is performed utilizing RNN structures over combined low-level verb-object representations (see Figure 2 (c)). In this experiment, either max-pooling or concatenation is applied to the verb softmax values from full-scale



756 and hand-scale verb models. Combined verb scores are concatenated with ob-  
757 ject softmax values, and final values are employed by the RNN action model.  
758 According to BiLSTM results, 46.47% and 45.18% accuracies are achieved  
759 for action recognition over split1 using max-pooling and concatenation, re-  
760 spectively. According to LSTM results, 44.10% and 44.36% accuracies are  
761 achieved for action recognition over split1 for max-pooling and concatena-  
762 tion, respectively. For both combination types, experimental results show  
763 that BiLSTM structure improves the accuracy compared to LSTM in Ta-  
764 ble 7. The best BiLSTM early fusion model with 46.47% mAcc is simply  
765 constructed using 1 layer with 100 cell size (see row 2 in Table 4). Please  
766 note that we do not conduct fc7 experiments for early fusion, since the di-  
767 mensions of fc7 verb features and object score vector  $\mathbf{s}^o$  are imbalanced for  
768 concatenation.

769  
770 **Recurrent verb-object late fusion.** In this fusion setting, action recog-  
771 nition is performed with a shallow neural network that takes the RNN en-  
772 codings of individual verb and object streams as the concatenated verb cat-  
773 egory score vector  $\mathbf{v}$  and object category score vector  $\mathbf{o}$  as inputs (see Fig-  
774 ure 2 (d)). Here, the verb vector  $\mathbf{v}$  is the output of the verb model *C3D*  
775 *full+hand(concatenation)-softmax-BiLSTM* (see Table 5) and the object vec-  
776 tor  $\mathbf{o}$  is the output of the object model *YOLO-s<sup>o</sup>-BiLSTM* (see Table 6). As  
777 given in Table 7, 45.45% accuracy is obtained over split1 with  $([\vec{\mathbf{v}}, \vec{\mathbf{o}}])$  (see  
778 row 3 in Table 4). This result is slightly lower than the early fusion strat-  
779 egy. Using the concatenated outputs of forward and backward functions,  
780  $[\vec{\mathbf{v}}, \overleftarrow{\mathbf{v}}, \vec{\mathbf{o}}, \overleftarrow{\mathbf{o}}]$ , the model performance later improves up to 49.01% mAcc.

781 Conducting experiments using fc7 layer as verb features, we obtain 46.62%  
782 accuracy using forward feature as input of the shallow network and 48.46%  
783 accuracy using forward and backward feature combination as input of the  
784 shallow network. While fc7 provides improvement over softmax feature with  
785 BiLSTM, it shows slightly worse performance for BiLSTM Bw+Fw. If we  
786 further extend the inputs with an object model using  $[\mathbf{s}^o, \mathbf{d}^o]$  representation  
787 (the other object models use  $\mathbf{s}^o$ ), the performance is 50.29%. Figure 5 shows  
788 a comparison chart reporting category based accuracies and results indicate  
789 that hand-object interaction through  $\mathbf{d}^o$  improves action recognition in 48  
790 categories with 1.83% again in accuracy. By analyzing the improved samples,  
791 we found that  $\mathbf{d}^o$  helps to correct the verb and action predictions for many  
792 video instances.  $\mathbf{d}^o$  models the layout of hands and objects with respect  
793 to each other within a frame and our RNN structure models the temporal

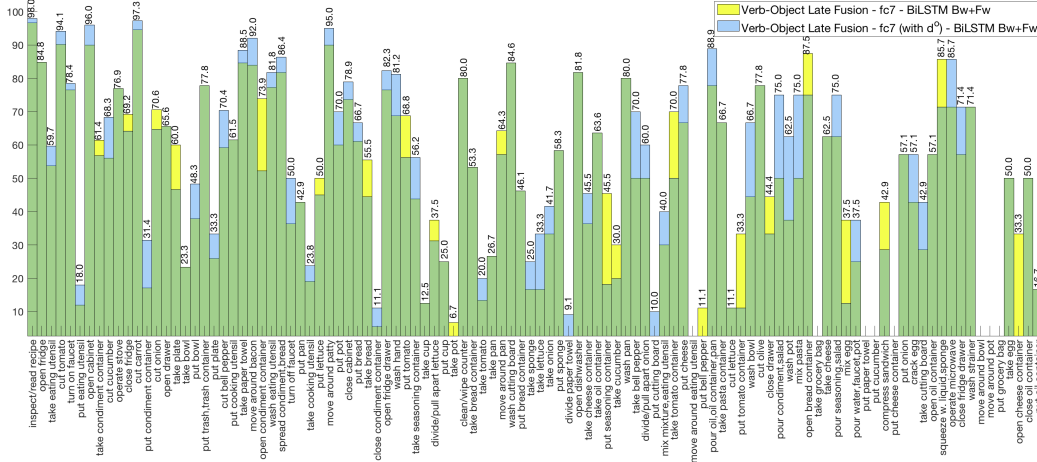


Figure 5: Comparison of Action Models *Verb-Object Late Fusion-fc7-BiLSTM Bw+Fw* with 48.46% accuracy and *Verb-Object Late Fusion-fc7 (with  $d^o$ )-BiLSTM Bw+Fw* with 50.29% accuracy (see Table 7). Maximum value per category over two models are shown (best viewed in color).

794 dynamics of the layout within frame and across frames. This helps to correct  
 795 verb and action predictions as shown in Figure 6.

796 These experiments show that (i) fc7 features, (ii) forward and backward  
 797 extension, and (iii) object model with spatial layout  $d^o$  improve the perfor-  
 798 mance significantly. Figure 7 shows same recognition results using the best  
 799 performed C3D model and how it predicts in challenging video cases. The  
 800 best shallow network with 50.29% mAcc consists of 3 fully connected layers  
 801 with 4096, 1024, 512 neurons and a 106-dimensional softmax output (see row  
 802 4 in Table 4).

803  
 804 **Recurrent action baseline.** Recurrent action baseline is also experimented  
 805 over C3D action model (see Figure 2 (f)). In this experiment, accuracy  
 806 reaches to 25.36% in Table 7. Result shows that the Recurrent Action Base-  
 807 line model is better than the Count-Based Action Baseline model due to  
 808 modelling of temporal dynamics in action videos (see Table 3), and the two-  
 809 stream recurrent fusion strategies that rely on either early or late fusion are  
 810 better than the baseline models.

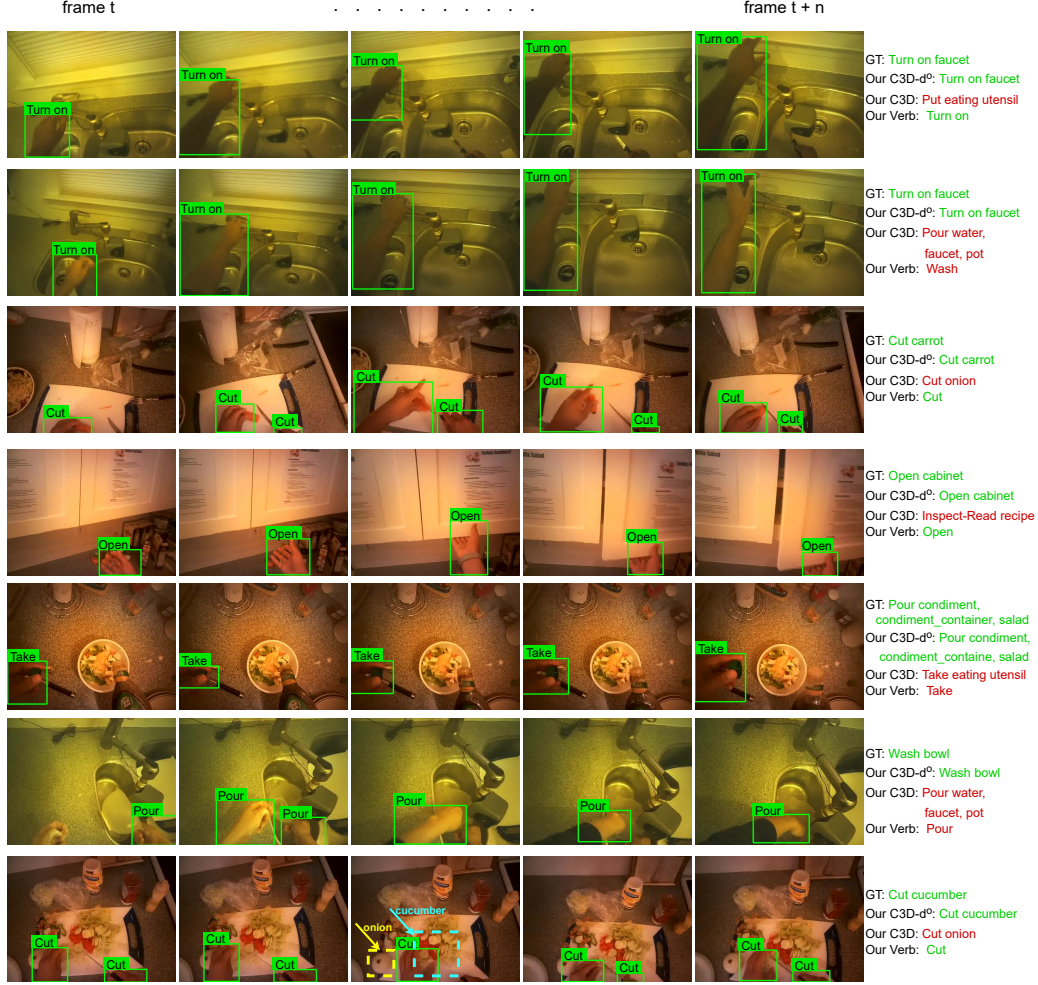


Figure 6: Visual comparison of action models for test samples from EGTEA Gaze+ dataset to evaluate the effect of spatial layout  $d^o$  (best viewed in color). For fair comparison, the late fusion strategies of *Our C3D-d°* (*Verb-Object Late Fusion-fc7(with  $d^o$ )-BiLSTM Bw+Fw*) and *Our C3D* (*Verb-Object Late Fusion-fc7-BiLSTM Bw+Fw*) (see Table 7) action models are considered. Both action models are constructed using same verb model. It can be easily appeared from the visual samples that action model with object model with distance scores  $d^o$  improves the action performances and corrects the verb and object predictions. For instance, for the video clip sample in the last row, although the verb action is predicted as *cut* and there is *onion* object in the background, *Our C3D* fails by predicting the video as *cut onion* action category. *Our C3D-d°* action model predicts action category of *cut cucumber* correctly since it takes into consideration the hand and *cucumber* object locations and interaction.



Figure 7: Visualization of same video samples to show the improvements of C3D action model for the failure cases of the verb model. The predictions of the best action model using *Our C3D-d° (Verb-Object Late Fusion-fc7(with d°)-BiLSTM Bw+Fw)*(see Table 7) and the predictions of its verb-RNN stream are illustrated. Even if the verb model predicts incorrectly, action model corrects the prediction of the verb stream. The verb predictions are confused due to the similar background of the video and the similarity of the hand movements. For instance, in the first row, the hand action is categorized as *cut* since the background and interacted objects are proper for that action although the ground truth verb action is *take*. Our C3D action model handles this failure cases of verb model and predicts the video correctly as *cut eating utensil* with correction of verb category.

Verb Base Model	Scale	Feature (pooled)	Verb-RNN	mAcc (%)
Pre-trained I3D	full	Mixed-5c (1x1x1x1024)	BiLSTM	66.03
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Pre-trained I3D	full	Mixed-5c (1x1x1x1024)	BiLSTM	71.62
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Pre-trained I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM	73.15
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Our I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM	<b>73.43</b>
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Our I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM+Attention	<b>74.81</b>
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		

Base Model	Scale	Feature (pooled)	Verb/Object-RNN	mAcc (%)
Pre-trained I3D	full	Mixed-5c (1x1x1x1024)	BiLSTM Bw+Fw	51.15
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Pre-trained I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM Bw+Fw	53.82
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Our I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM Bw+Fw	<b>54.56</b>
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		
Our I3D	full	Mixed-5c (3x1x1x1024)	BiLSTM (Attention) Bw+Fw	<b>56.07</b>
	hand	MaxPool3d-5a-2x2 (1x1x1x832)		

Table 8: I3D verb and action recognition results on split1 videos with recurrent-based models. The action models are based on recurrent verb-object late fusion strategy, and trained using backward and forward function outputs. The object stream used in action models is based on model  $YOLO-s^o + d^o-BiLSTM$ . The reported results as the last verb and action models are trained with attention module. Attention module is applied both on verb and object models, but we obtain no improvement for object recognition.

#### 811 4.4.4. Other 3D ConvNet Architectures: I3D

812 The base encoding models used for the feature extraction may improve the  
813 recognition rates significantly. Following this, we extend our experiments by  
814 using other 3D ConvNet architecture called I3D RGB [40] to categorize clips  
815 (following the original setting we use 25-frame clips) as the base model for the  
816 verb stream and we conduct two experiments. First, we use the pre-trained  
817 I3D model that is trained on the Kinetics dataset with 400 action categories  
818 <sup>2</sup>. Later, we train a shallow neural network model that fine-tunes over I3D  
819 intermediate features using our verb and action annotations. In both cases,  
820 the base model encodes videos as verb matrices and we train our verb and  
821 action models using recurrent models on top of these matrices as before. We  
822 report the verb and the action recognition accuracies, respectively, in Table  
823 8. Please note that we conduct I3D experiments with the setting that is best  
824 performed for C3D, therefore we investigate the Recurrent Verb-Object Late  
825 Fusion strategy that performs best for C3D (given in Table 7) and the object

<sup>2</sup><https://github.com/deepmind/kinetics-i3d>



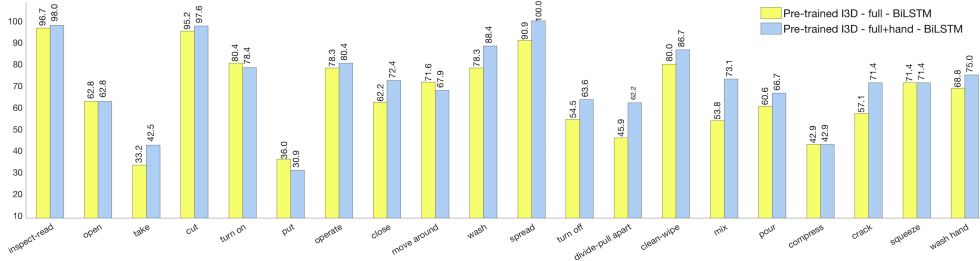


Figure 8: Comparison of multiple scale I3D verb models over 20 verb categories of the EGTEA Gaze+ dataset. Yellow bars show the recognition results of verb model *Pre-trained I3D-full-BiLSTM* and blue shows the results of verb model *Pre-trained I3D-full+hand-BiLSTM* (see in Table 8). Multiple scales using full and hand scales improve recognition with 5.59% gain.

stream is object model  $YOLO-s^o+d^o-BiLSTM$ . Figure 9 shows predictions over a set of test samples using I3D and C3D models. When we analyze the samples, we observe I3D resolves some challenging cases and improves recognition performance over C3D. I3D is a dense network compared to C3D and intermediate layers consists of inception modules. Using a more advance base model for verb stream, we improve recognition performance by 4.27%.

For I3D experiments, two intermediate layers from the I3D model are selected for encoding the full-scale and the hand-scale verb information, respectively. For the full-scale encoding used to fetch coarser details in clips, we pick the outputs of  $3 \times 7 \times 7 \times 1024$ -dimensional Mixed-5c layer. On the other hand, for the hand-scale encoding corresponding to finer details on hands, we pick the outputs of an earlier layer,  $3 \times 14 \times 14 \times 832$ -dimensional MaxPool3d-5a-2x2 layer for hand-volumes. We apply 20% enlargement on detected hand regions before extracting features of hand-volumes. Finally, each clip is represented as the concatenation of these features. If there are multiple hands, we apply max-pooling on features of hand-volumes.

Using pre-trained I3D model, we first examine how the multiple scales help in performance and compare full-scale (66.03% mAcc) vs. combination of full-scale and hand-scale models (71.62% mAcc). We also show the details for all verb categories in Figure 8. There is a significant improvement over 18 verb categories with the addition of fine-grained details through hand-scale model. This is what we expect from the combination of fine and coarse scale features. Then, we experiment on two different feature setting. We obtain 71.62% mAcc and 73.15% mAcc for verb recognition (the BiLSTM

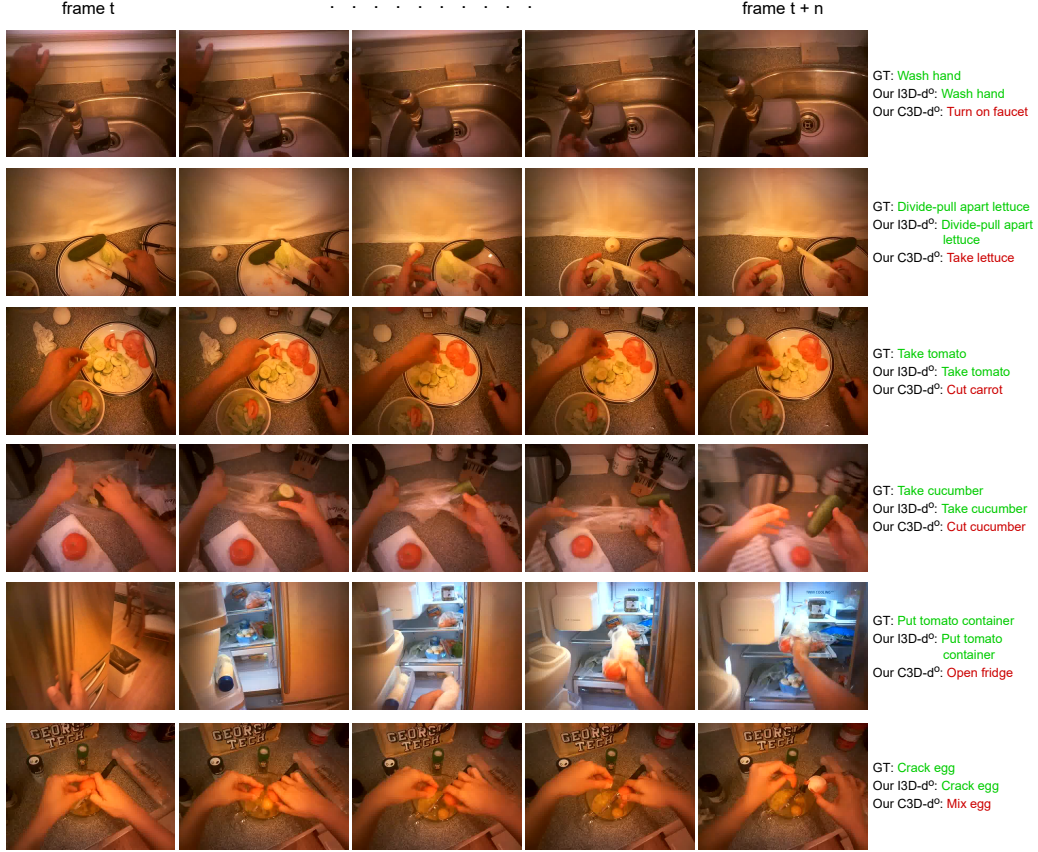


Figure 9: Visual comparison of predictions of I3D and C3D based action models on test samples from EGTEA Gaze+ dataset (best viewed in color). For fair comparison, the same late fusion settings of C3D- $d^o$  (*Verb-Object Late Fusion- $fc7(with\ d^o)$ -BiLSTM Bw+Fw*) from Table 7 and I3D- $d^o$  (*Our I3D-full+hand-BiLSTM Bw+Fw*) from Table 8 action models are considered. Both models are constructed using the same object model. These challenging videos are predicted true by I3D action model while predicted false by C3D action model. The illustration of successive frames of videos shows the challenging videos are truly predicted by I3D while missed by C3D. For the video sample in the first row, while the *turn-on faucet* action is performed at the beginning of the video, I3D based action model predicts the main video action correctly. The same challenge is also valid for the second and fifth row video samples. In the third row video sample, both object and verb categories are wrong in C3D. Although the background is proper for *cut* verb action, I3D overcomes this challenging situation.

verb models contain 1 layer with 728 cells). Moreover, we obtain 51.15% mAcc and 53.82% mAcc in action recognition (the action models contain 2 1024-dimensional fully connected layers with dropout rate of 0.5). In the first feature setting, the full-scale features and the hand-scale features are pooled both in spatial and temporal dimensions, and then concatenated. In the second one, the full-scale features are pooled just in the spatial dimension and concatenated with the pooled hand-volume features. Results show that higher dimensional representations encode the data better. Please note that further experiments can be conducted on other network layers with various pooling settings to improve the recognition performance. Keeping the features as in the original dimension is good to encode spatial and temporal information, but here we prefer to apply pooling over spatial domain to decrease feature dimensions.

Our fine-tuned model is trained over extracted I3D features. Unlike two separate models introduced on C3D for full-scale verb and hand-scale verb features respectively, here we train a single model on concatenated I3D full-scale and hand-scale features. Representing each clip with a concatenated feature vector, we train a verb model over ground truth action clips. This model contains 2 1024-dimensional fully connected layers with dropout rate of 0.5. Using this verb model, we represent each verb category score vector  $\mathbf{v}$  using the 1024-dimensional second fully connected layer output. Later,  $\mathbf{v}$  and  $\mathbf{o}$  are concatenated and fed into the shallow network. This shallow network similarly contains 2 1024-dimensional fully connected layers with dropout rate of 0.5 (see row 5 in Table 4). The verb and the action models result in recognition accuracies of 73.43% and 54.56%, respectively. Our supervised setting slightly performs better than the pre-trained models, since it provides fine-tuning over EGTEA Gaze+ dataset.

#### 4.4.5. Recurrent Models with Attention

As proposed in Section 3.3, we improve fusion strategies with a self-attention module and we propose a model called Recurrent Verb-Object Attention Late Fusion strategy. The attention module can be easily applied to early fusion strategy that reduces verb and object streams into a single pathway in earlier stages, but we work on late fusion since it is superior in recognition performance (see Table 7). Please note that we conduct experiments for the best performing setting, therefore attention module is applied for late fusion strategy and we use our I3D models. The verb and object BiLSTM models with attention module are constructed using 1 layer with 512



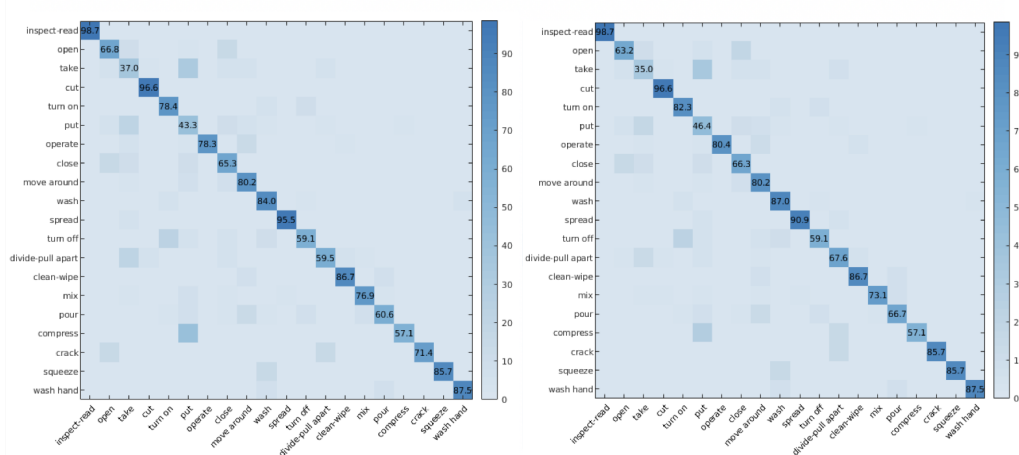


Figure 10: Confusion matrices reporting the performances of verb models *Our I3D-full+hand-BiLSTM* (left) and *Our I3D-full+hand-BiLSTM+Attention* (right) with 73.43% and 74.81% mAcc, respectively, over 20 verb categories of the EGTEA Gaze+ dataset (see Table 8).

and 256 cell sizes, respectively. Finally, the action model using Verb-Object Attention Late Fusion setting contains two 2048-dimensional fully connected layers with dropout rate of 0.5 (see row 6 in Table 4).

In Verb-Object Attention Late Fusion setting, the verb, object and action model results are 74.81%, 70.83% and 56.07%, respectively. According to the results, the recurrent verb and action models with attention block outperform the non-attention models as seen in Table 8. The verb model results in 1.38% gain and action model results in 1.51% gain in recognition accuracies. However, we observe no improvement on object model. For comparison, more details can be found in Figure 10 and in Figure 11. We also show some failure cases in Figure 12. The failures are caused by the similarity of the object categories and hand movements.

## 5. Comparison and Discussion

This section presents the comparison of the proposed model on the EGTEA Gaze+ dataset for action recognition, as well as a detailed discussion.

Table 9 reports results of our proposed method with the state-of-the-art models and it shows that our performance is comparable with the state-of-the-art. Baseline models I3D RGB, I3D Joint, I3D+Gaze returns 47.26%, 49.79%, and 51.21% accuracies, respectively. Results show that I3D is a

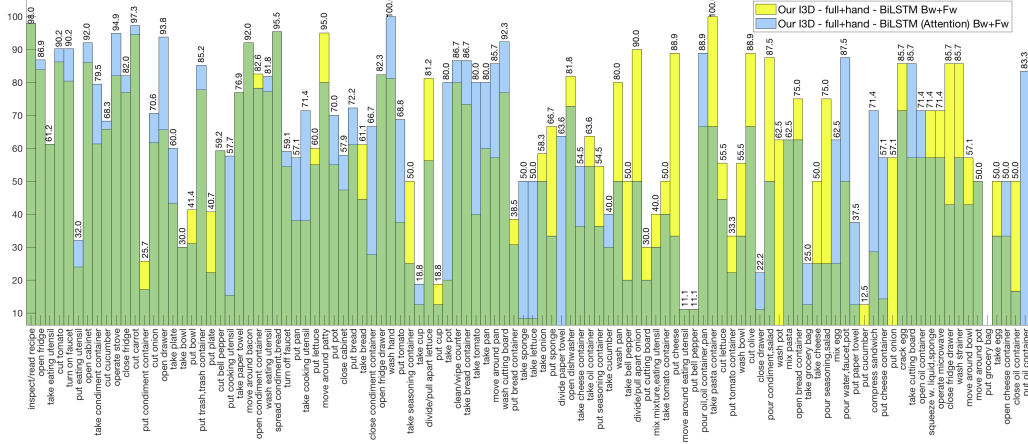


Figure 11: Comparison of action models *Our I3D-full+hand-BiLSTM Bw+Fw* and *Our I3D-full+hand-BiLSTM(Attention) Bw+Fw* with 54.56% and 56.07% mAcc over 106 action categories of the EGTEA Gaze+ dataset (best viewed in color).

Action Models	mAcc (%)
EgoIDT + Gaze [30]	46.50
I3D RGB [13]	47.26
I3D Joint [13]	49.79
I3D+Gaze [13]	51.21
Li et al. [13]	53.30
MCN [14]	55.63
<b>Ours</b>	<b>56.07</b>
LSTA-RGB [45]	57.94
RU [46]	60.20
LSTA [45]	61.86

Table 9: The comparison with state-of-the-art action recognition models on the first split of EGTEA Gaze+ dataset with accuracy in mAcc.

powerful feature, and I3D Joint with joint modelling of RGB and Flow features improves the recognition significantly. Moreover, I3D+Gaze has the highest accuracy among three, since gaze is an important clue for egocentric videos. Our models rely only on RGB modality, and they are better than the I3D based models and Li et al. [13]. This means the performance gap can be increased with integration of other modalities into our pipeline. Integration is simple, where each stream of our verb-object decomposition model can be further extended with two-stream approaches to add Flow or Gaze

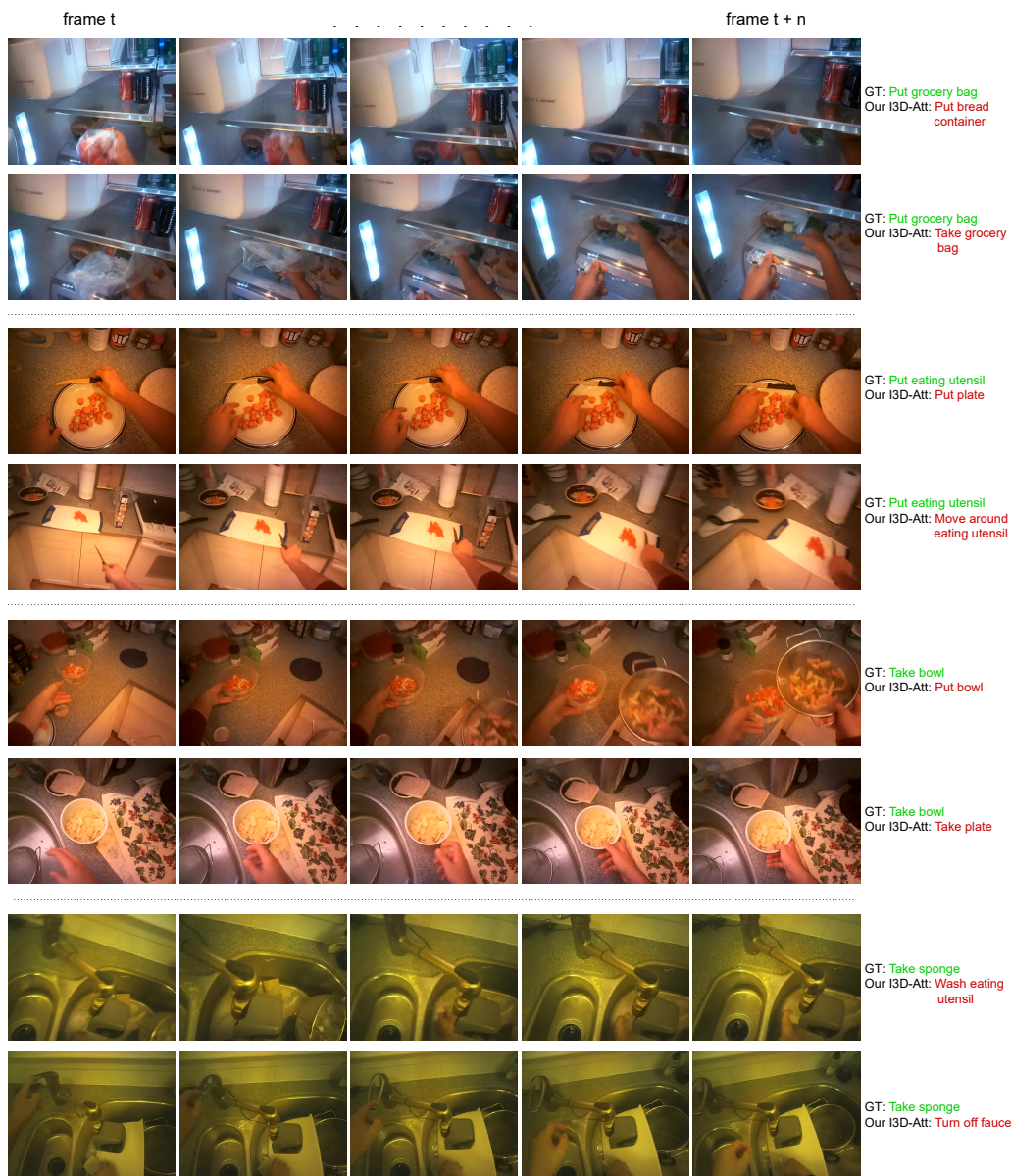


Figure 12: Failure cases of *I3D-Att* (*Our I3D-full+hand-BiLSTM(Attention) Bw+Fw*, see Table 8) action model on a set of test samples from split1.

914 modalities.

915 We report lower results than the recently published studies the RU [46]  
916 and the LSTA [45]. In [46], the Rolling-Unrolling LSTM (RU) processes  
917 appearance from RGB frames, motion from optical flow, as well as object  
918 features. These modalities are fused using an attention mechanism. Simi-  
919 larly, LSTA [45] is a two-stream model; one stream for encoding appearance  
920 information from RGB frames and the second stream for encoding motion  
921 information from optical flow. Both models use flow as a low-level feature.  
922 Following a similar discussion, we state that our model focuses on a single  
923 modality with RGB features. The standalone performance of LSTA RGB  
924 stream, LSTA-RGB, [45] has 57.94% mAcc and this is slightly higher than  
925 our best model with 56.07%. Please note that we use the term two-stream  
926 to indicate the verb-object decomposition of action model, but other mod-  
927 els use the term of two-stream architecture following the work [32] where  
928 each stream models a different low-level modality, namely RGB and Flow.  
929 This means that the performance of object and verb streams can be further  
930 improved with other modalities.

931 Investigating fine-grained recognition in first-person view, our aim is on  
932 how the recognition rates can be improved within the model just using RGB  
933 modality. We focus on hands, their actions in multiple scales and their  
934 interactions with other objects. Verb stream is modelled using RGB with a  
935 single modality and 3D convolutional neural network models are investigated  
936 for modelling multiple scales of hand regions. Similarly, object-stream is  
937 modelled using RGB frames and spatial layout features.

938 On the other hand, our proposed model is based on action decomposition  
939 with two semantically meaningful components, verb and object. In this work,  
940 we show that we achieve comparable results with the state-of-the-art models.  
941 Even if we have slightly lower performance than some of the approaches, our  
942 model has many architectural advantages over conventional action recogni-  
943 tion models. First, decomposition is good for zero-shot learning as proposed  
944 in [7]. Second, in large scale datasets, the number of video instance can vary  
945 for each category, and while some categories have a large amount of train-  
946 ing samples, some other categories have very few training samples. Through  
947 decomposition based models, we have less number of categories with more  
948 samples to train each component of neural network architectures, and this  
949 helps to solve the problem related to dataset imbalance. Moreover, models  
950 based on action decomposition are architecturally more flexible for extending  
951 the model later for more categories. For example, fixing the object compo-

952 nent with trained model, an addition of a new verb category will cost to  
953 fine-tune the verb component and the simple fusion model (observing late  
954 fusion outperforms early fusion).

## 955 6. Conclusion

956 We have developed compositional model including two complementary  
957 steps, verb and object, to perform action recognition in first-person videos.  
958 Late and early fusion strategies, based on recurrent neural network struc-  
959 tures pools verb models and object model to recognize the video at action  
960 level. Experimental results show that decomposing actions model into verb  
961 and object using recurrent neural networks significantly improves the per-  
962 formance compared to the baseline action model for large number of action  
963 classes. Hand information is an important clue to determine the action in  
964 the first-person vision. We have shown that spatial-temporal modelling of  
965 hand regions improves both verb and action recognition performances.

## 966 References

- 967 [1] R. Megret, D. Szolgay, J. Benois-Pineau, P. Joly, J. Piquier, J.-F. Dar-  
968 tiges, C. Helmer, Wearable video monitoring of people with age de-  
969 mentia: Video indexing at the service of helthcare, in: International  
970 Workshop on Content-based Multimedia Indexing, 2008, pp. 101–108.
- 971 [2] G. Meditskos, P.-M. Plans, T. G. Stavropoulos, J. Benois-Pineau,  
972 V. Buso, I. Kompatsiaris, Multi-modal activity recognition from ego-  
973 centric vision, semantic enrichment and lifelogging applications for the  
974 care of dementia, *Journal of Visual Communication and Image Repre-*  
975 *sentation* 51 (2018) 169–190.
- 976 [3] N. Das, E. Ohn-Bar, M. M. Trivedi, On performance evaluation of  
977 driver hand detection algorithms: Challenges, dataset, and metrics, in:  
978 International Conference on Intelligent Transportation Systems, 2015,  
979 pp. 2953–2958.
- 980 [4] J. Lee, M. S. Ryoo, Learning robot activities from first-person human  
981 videos using convolutional future regression, in: IEEE Conference on  
982 Computer Vision and Pattern Recognition Workshops, 2017, pp. 1–2.

- 983 [5] H. Rhodin, C. Richardt, D. Casas, E. Insafutdinov, M. Shafiei, H.-P.  
984 Seidel, B. Schiele, C. Theobalt, Egocap: egocentric marker-less mo-  
985 tion capture with two fisheye cameras, *ACM Transactions on Graphics*  
986 (TOG) 35 (2016) 162.
- 987 [6] M. Ma, H. Fan, K. M. Kitani, Going deeper into first-person activity  
988 recognition, in: *IEEE Conference on Computer Vision and Pattern*  
989 *Recognition*, 2016, pp. 1894–1903.
- 990 [7] Y. C. Zhang, Y. Li, J. M. Rehg, First-person action decomposition  
991 and zero-shot learning, in: *IEEE Winter Conference on Applications of*  
992 *Computer Vision*, IEEE, 2017, pp. 121–129.
- 993 [8] H. Wang, C. Schmid, Action recognition with improved trajectories, in:  
994 *Proceedings of the IEEE international conference on computer vision*,  
995 2013, pp. 3551–3558.
- 996 [9] A. Fathi, X. Ren, J. M. Rehg, Learning to recognize objects in egocen-  
997 tric activities, in: *IEEE Conference on Computer Vision and Pattern*  
998 *Recognition*, 2011, pp. 3281–3288.
- 999 [10] M. Cai, F. Lu, Y. Gao, Desktop action recognition from first-person  
1000 point-of-view, *IEEE Transactions on Cybernetics* 49 (2018) 1616–1628.
- 1001 [11] T. Ishihara, K. M. Kitani, W.-C. Ma, H. Takagi, C. Asakawa, Recog-  
1002 nizing hand-object interactions in wearable camera videos, in: *IEEE*  
1003 *International Conference on Image Processing*, 2015.
- 1004 [12] J. Kumar, Q. Li, S. Kyal, E. A. Bernal, R. Bala, On-the-fly hand de-  
1005 tection training with application in egocentric action recognition, in:  
1006 *IEEE Conference on Computer Vision and Pattern Recognition Work-*  
1007 *shops*, 2015, pp. 18–27.
- 1008 [13] Y. Li, M. Liu, J. M. Rehg, In the eye of beholder: Joint learning of  
1009 gaze and actions in first person video, in: *Proceedings of the European*  
1010 *Conference on Computer Vision*, 2018, pp. 619–635.
- 1011 [14] Y. Huang, Z. Li, M. Cai, Y. Sato, Mutual context network for jointly  
1012 estimating egocentric gaze and actions, *arXiv preprint arXiv:1901.01874*  
1013 (2019).

- 1014 [15] E. H. Spriggs, F. De La Torre, M. Hebert, Temporal segmentation and  
1015 activity classification from first-person sensing, in: IEEE Conference on  
1016 Computer Vision and Pattern Recognition Workshops, 2009, pp. 17–24.
- 1017 [16] A. Fathi, A. Farhadi, J. M. Rehg, Understanding egocentric activities,  
1018 in: IEEE International Conference on Computer Vision, 2011, pp. 407–  
1019 414.
- 1020 [17] H. Pirsiavash, D. Ramanan, Detecting activities of daily living in first-  
1021 person camera views, in: IEEE Conference on Computer Vision and  
1022 Pattern Recognition, IEEE, 2012, pp. 2847–2854.
- 1023 [18] M. S. Ryoo, B. Rothrock, L. Matthies, Pooled motion features for first-  
1024 person videos, in: IEEE Conference on Computer Vision and Pattern  
1025 Recognition, 2015, pp. 896–904.
- 1026 [19] Y. Poley, A. Ephrat, S. Peleg, C. Arora, Compact cnn for indexing  
1027 egocentric videos, in: IEEE Winter Conference on Applications of Com-  
1028 puter Vision, IEEE, 2016, pp. 1–9.
- 1029 [20] Y. Zhou, B. Ni, R. Hong, X. Yang, Q. Tian, Cascaded interactional  
1030 targeting network for egocentric video analysis, in: IEEE Conference on  
1031 Computer Vision and Pattern Recognition, 2016, pp. 1904–1913.
- 1032 [21] S. Singh, C. Arora, C. Jawahar, Trajectory aligned features for first  
1033 person action recognition, Pattern Recognition 62 (2017) 45–55.
- 1034 [22] R. E. Schapire, Explaining adaboost, in: Empirical inference, Springer,  
1035 2013, pp. 37–52.
- 1036 [23] N. Dalal, B. Triggs, Histograms of oriented gradients for human detec-  
1037 tion, in: IEEE Conference on Computer Vision and Pattern Recognition,  
1038 2005.
- 1039 [24] A. Cartas, P. Radeva, M. Dimiccoli, Contextually driven first-person  
1040 action recognition from videos, in: presentation at EPIC@ ICCV2017  
1041 workshop, 2017, p. 8.
- 1042 [25] J. Pont-Tuset, P. Arbelaez, J. T. Barron, F. Marques, J. Malik, Mul-  
1043 tiscscale combinatorial grouping for image segmentation and object pro-  
1044 posal generation, IEEE Transactions on Pattern Analysis and Machine  
1045 Intelligence 39 (2017) 128–140.

- 1046 [26] G. Gkioxari, R. B. Girshick, J. Malik, Contextual action recognition  
1047 with r\*cnn, in: IEEE International Conference on Computer Vision,  
1048 2015, pp. 1080–1088.
- 1049 [27] A. Fathi, Y. Li, J. M. Rehg, Learning to recognize daily actions using  
1050 gaze, in: European Conference on Computer Vision (ECCV), 2012, pp.  
1051 314–327.
- 1052 [28] Y. Li, A. Fathi, J. M. Rehg, Learning to predict gaze in egocentric  
1053 video, in: IEEE International Conference on Computer Vision, 2013,  
1054 pp. 3216–3223.
- 1055 [29] K. M. Kitani, T. Okabe, Y. Sato, A. Sugimoto, Fast unsupervised ego-  
1056 action learning for first-person sports videos, in: IEEE Conference on  
1057 Computer Vision and Pattern Recognition, 2011, pp. 3241–3248.
- 1058 [30] Y. Li, Z. Ye, J. M. Rehg, Delving into egocentric actions, in: IEEE  
1059 Conference on Computer Vision and Pattern Recognition, 2015, pp.  
1060 287–295.
- 1061 [31] H. Wang, A. Kläser, C. Schmid, L. Cheng-Lin, Action recognition by  
1062 dense trajectories, in: IEEE Conference on Computer Vision and Pat-  
1063 tern Recognition, 2011, pp. 3169–3176.
- 1064 [32] K. Simonyan, A. Zisserman, Two-stream convolutional networks for  
1065 action recognition in videos, in: Advances in Neural Information Pro-  
1066 cessing Systems, 2014, pp. 568–576.
- 1067 [33] Y. Tang, Y. Tian, J. Lu, J. Feng, J. Zhou, Action recognition in rgb-d  
1068 egocentric videos, in: International Conference on Image Processing,  
1069 2017, pp. 3410–3414.
- 1070 [34] M. Hahn, N. Ruiz, J.-B. Alayrac, I. Laptev, J. M. Rehg, Learning  
1071 to localize and align fine-grained actions to sparse instructions, arXiv  
1072 preprint arXiv:1809.08381 (2018).
- 1073 [35] G. Kapidis, R. Poppe, E. van Dam, L. P. Noldus, R. C. Veltkamp, Ego-  
1074 centric hand track and object-based human action recognition, arXiv  
1075 preprint arXiv:1905.00742 (2019).



- 1076 [36] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image  
1077 recognition, in: IEEE Conference on Computer Vision and Pattern  
1078 Recognition, 2016, pp. 770–778.
- 1079 [37] D. Damen, H. Doughty, G. Maria Farinella, S. Fidler, A. Furnari,  
1080 E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, et al., Scal-  
1081 ing egocentric vision: The epic-kitchens dataset, in: Proceedings of the  
1082 European Conference on Computer Vision, 2018, pp. 720–736.
- 1083 [38] J. Redmon, A. Farhadi, Yolov3: An incremental improvement, arXiv  
1084 preprint arXiv:1804.02767 (2018).
- 1085 [39] D. Tran, L. Bourdev, R. Fergus, L. Torresani, M. Paluri, Learning  
1086 spatiotemporal features with 3d convolutional networks, in: IEEE In-  
1087 ternational Conference on Computer Vision, 2015, pp. 4489–4497.
- 1088 [40] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model  
1089 and the kinetics dataset, in: IEEE Conference on Computer Vision and  
1090 Pattern Recognition, 2017, pp. 6299–6308.
- 1091 [41] J. Redmon, A. Farhadi, Yolo9000: better, faster, stronger, in: IEEE  
1092 Conference on Computer Vision and Pattern Recognition, 2017, pp.  
1093 7263–7271.
- 1094 [42] L. Shen, S. Yeung, J. Hoffman, G. Mori, L. Fei-Fei, Scaling human-  
1095 object interaction recognition through zero-shot learning, in: IEEE  
1096 Winter Conference on Applications of Computer Vision, 2018, pp. 1568–  
1097 1576.
- 1098 [43] Z. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Zhou, Y. Ben-  
1099 gio, A structured self-attentive sentence embedding, arXiv preprint  
1100 arXiv:1703.03130 (2017).
- 1101 [44] Georgia tech egocentric activity datasets,  
1102 <http://www.cbi.gatech.edu/fpv/>, ??? Accessed: 2019-06-14.
- 1103 [45] S. Sudhakaran, S. Escalera, O. Lanz, Lsta: Long short-term attention  
1104 for egocentric action recognition, in: IEEE Conference on Computer  
1105 Vision and Pattern Recognition, 2019, pp. 9954–9963.

- 1106 [46] A. Furnari, G. M. Farinella, What would you expect? anticipating  
1107 egocentric actions with rolling-unrolling lstms and modality attention,  
1108 in: IEEE International Conference on Computer Vision, 2019, pp. 6252–  
1109 6261.