

# Hunting Group Clues with Transformers for Social Group Activity Recognition

Masato Tamura<sup>✉</sup>, Rahul Vishwakarma<sup>✉</sup>, and Ravigopal Vennelakanti

Hitachi America, Ltd.

masato.tamura.sf@hitachi.com,

{rahul.vishwakarma,ravigopal.vennelakanti}@hal.hitachi.com

**Abstract.** This paper presents a novel framework for social group activity recognition. As an expanded task of group activity recognition, social group activity recognition requires recognizing multiple sub-group activities and identifying group members. Most existing methods tackle both tasks by refining region features and then summarizing them into activity features. Such heuristic feature design renders the effectiveness of features susceptible to incomplete person localization and disregards the importance of scene contexts. Furthermore, region features are sub-optimal to identify group members because the features may be dominated by those of people in the regions and have different semantics. To overcome these drawbacks, we propose to leverage attention modules in transformers to generate effective social group features. Our method is designed in such a way that the attention modules identify and then aggregate features relevant to social group activities, generating an effective feature for each social group. Group member information is embedded into the features and thus accessed by feed-forward networks. The outputs of feed-forward networks represent groups so concisely that group members can be identified with simple Hungarian matching between groups and individuals. Experimental results show that our method outperforms state-of-the-art methods on the Volleyball and Collective Activity datasets.

**Keywords:** social group activity recognition, group activity recognition, social scene understanding, attention mechanism, transformer

## 1 Introduction

Social group activity recognition is a task of recognizing multiple sub-group activities and identifying group members in a scene. This task is derived from group activity recognition, which needs to recognize only one group activity in a scene. Both tasks have gained tremendous attention in recent years for potential applications such as sports video analysis, crowd behavior analysis, and social scene understanding [1–5, 12–14, 16–18, 21, 23–27, 32, 33, 36, 40–43]. In the context of these tasks, the term “action” denotes an atomic movement of a single person, and the term “activity” refers to a more complex relation of movements performed by a group of people. Although our framework can recognize both actions and activities, we focus on group activities.

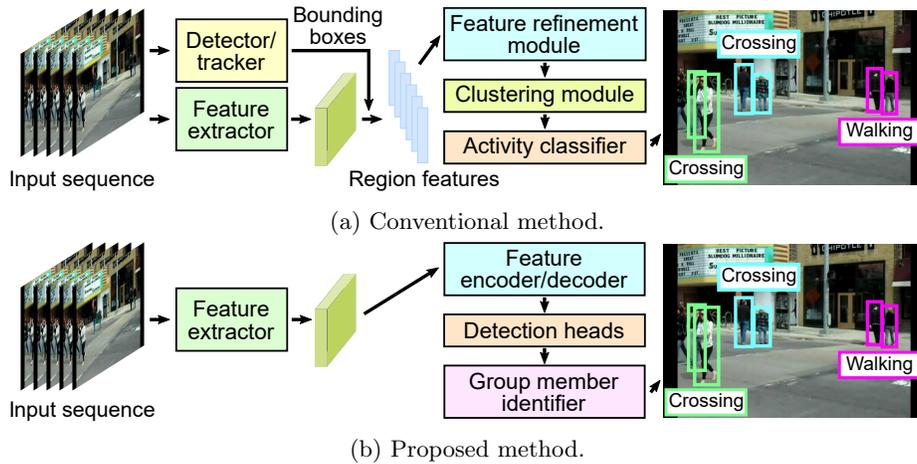


Fig. 1: Overviews of conventional and proposed social group activity recognition methods. The labels in the right image show predicted social group activities.

Most existing methods decompose the recognition process into two independent parts; person localization and activity recognition (See Fig. 1a) [5, 12–14, 16, 18, 21, 26, 32, 33, 36, 40–43]. Person localization identifies regions where people are observed in a scene with bounding boxes. These boxes are used to extract region features from feature maps. The region features are further refined to encode spatio-temporal relations with feature refinement modules such as recurrent neural networks (RNNs) [8, 15], graph neural networks (GNNs) [20, 39], and transformers [38]. The refined features are summarized for the purpose of activity recognition.

While these methods have demonstrated significant improvement, they have several drawbacks attributed to the heuristic nature of feature design. Since region features are extracted from bounding box regions in feature maps, the effectiveness of the features is affected by the localization performance. Most existing methods ignore this effect and evaluate their performances with region features of ground truth boxes. However, several works [5, 13, 33, 41] show that the recognition performance is slightly degraded when using predicted boxes instead of ground truth boxes. Moreover, substantial scene contexts are discarded by using region features because they are typically dominated by features of the people in the boxes. Scene contexts such as object positions and background situations are sometimes crucial to recognize group activities. For instance, positions of sports balls are informative to recognize group activities in sports games. These features should be leveraged to enhance recognition performance.

Another challenge specific to social group activity recognition is that utilizing region features is sub-optimal to identify group members. Ehsanpour *et al.* [13] use region features as node features of graph attention networks (GATs) [39] and train them to output adjacency matrices that have low probabilities for people

in different groups and high probabilities for those in the same groups. During inference, spectral clustering [31] is applied to the adjacency matrices to divide people into groups. Because adjacency matrices reflect semantic similarities of node features, this method may not work if region features of people in the same group have different semantics such as doing different actions.

To address these challenges, we propose a novel social group activity recognition method that can be applied to both social group activity recognition and group activity recognition. We leverage a transformer-based object detection framework [6, 45] to obviate the need for the heuristic feature design in existing methods (See Fig. 1b). Attention modules in transformers play crucial roles in our method. We design our method in such a way that the attention modules identify and then aggregate features relevant to social group activities, generating an effective feature for each social group. Because activity and group member information is embedded into the generated features, the information can be accessed by feed-forward networks (FFNs) in the detection heads. The outputs of the detection heads are designed so concisely that group member identification can be performed with simple Hungarian matching between groups and individuals. This identification method differs from Ehsanpour *et al.*'s method [13] in that their method relies on individuals' features to divide people into groups, while our method generates features that are embedded with clues for grouping people, enabling effective group identification.

To summarize, our contributions are three-fold:

- We propose a novel social group activity recognition method that leverages the attention modules in transformers to generate effective social group features. The group member information extracted from the features is designed to be concise and can be used to identify group members with a simple matching process.
- Our method achieves better or competitive performance to state-of-the-art methods on both group activity recognition and social group activity recognition in two challenging benchmarks.
- We perform comprehensive analyses to reveal how our method works with activities under various conditions.

## 2 Related Works

### 2.1 Group Activity Recognition

Deep-neural-network-based methods have become dominant in group activity recognition due to the learning capability of the networks. Ibrahim *et al.* [18] proposed an RNN-based method that uses convolutional neural networks to extract features of person bounding box regions and long short-term memories to refine region features. This architecture captures the temporal dynamics of each person between frames and spatial dynamics of people in a scene. After their work, several RNN-based methods were proposed [5, 21, 33, 36, 40].

GNNs are also utilized to model the spatio-temporal context and relationships of people in a scene. Wu *et al.* [41] used graph convolutional networks (GCNs) [20] to capture spatio-temporal relations of people’s appearances and positions between frames. Ehsanpour *et al.* [13] adopted GATs [39] to learn underlying interactions and divide people into social groups with adjacency matrices. Hu *et al.* [16] utilized both RNNs and GNNs with reinforcement learning to refine features. Yuan *et al.* [42] used person-specific dynamic graphs that dynamically change connections of GNNs for each node.

With the rapid application of transformers [38] to vision problems, several works introduced transformers into group activity recognition. Gavriluk *et al.* [14] used transformer encoders to refine region features. Li *et al.* [26] proposed spatial-temporal transformers that can encode spatio-temporal dependence and decode the group activity information. Zhou *et al.* [43] proposed multi-scale spatio-temporal stacked transformers for compositional understanding and relational reasoning in group activities.

Our method differs from existing methods in that they rely on region features, while our method generates social group features with the attention modules in transformers, resulting in improving the performance.

## 2.2 Detection Transformer

Carion *et al.* [6] proposed a transformer-based object detector called DETR, which regards object detection as a set prediction and achieves end-to-end object detection. One significant difference between conventional object detectors and DETR is that conventional ones need heuristic detection points whose features are used to predict object classes and bounding boxes, while DETR obviates such heuristic components by letting queries in transformer decoders aggregate features for their target objects with the attention mechanisms. DETR shows competitive performance compared with conventional state-of-the-art detectors even without such heuristic components.

To further improve the performance of DETR, several methods have been proposed [11, 37, 45]. Zhu *et al.* [45] proposed Deformable DETR that replaces standard transformers with deformable ones. Deformable attention modules in the transformers combine a sparse sampling of deformable convolution [10] and dynamic weighting of standard attention modules, which significantly reduces the computational complexity of the attention weight calculation. This reduction allows Deformable DETR to use multi-scale feature maps from backbone networks. To leverage non-heuristic designs and multi-scale feature maps, we use deformable transformers to generate social group features.

## 3 Proposed Method

We leverage a deformable-transformer-based object detection framework [45] to recognize multiple group activities and identify group members without the heuristic feature design. We first explain the overall architecture in Sec. 3.1

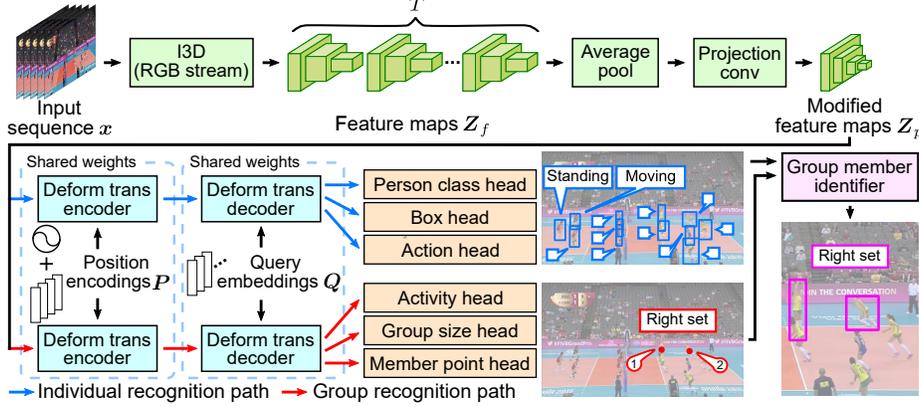


Fig. 2: Overall architecture of the proposed method.

and show how we set up the framework for social group activity recognition. In Sec. 3.2, we describe the loss function used in training. Finally, we explain the group member identification performed during inference in Sec. 3.3. Due to the limited space, we omit the details of deformable transformers and encourage readers to refer to the paper of Deformable DETR [45] for more details.

### 3.1 Overall Architecture

Figure 2 shows the overall architecture of the proposed method. Given a frame sequence  $\mathbf{x} \in \mathbb{R}^{3 \times T \times H \times W}$ , a feature extractor extracts a set of multi-scale feature maps  $\mathbf{Z}_f = \{\mathbf{z}_i^{(f)} \mid \mathbf{z}_i^{(f)} \in \mathbb{R}^{D_i \times T \times H'_i \times W'_i}\}_{i=1}^{L_f}$ , where  $T$  is the length of the sequence,  $H$  and  $W$  are the height and width of the frame,  $H'_i$  and  $W'_i$  are those of the output feature maps,  $D_i$  is the number of channels, and  $L_f$  is the number of scales. We adopt the inflated 3D (I3D) network [7] as a feature extractor to embed local spatio-temporal context into feature maps. Note that we use only the RGB stream of I3D because group members are identified by their positions, which cannot be predicted with the optical flow stream. To reduce the computational costs of transformers, each feature map  $\mathbf{z}_i^{(f)}$  is mean-pooled over the temporal dimension and input to a projection convolution layer that reduces the channel dimension from  $D_i$  to  $D_p$ . One additional projection convolution layer with a kernel size of  $3 \times 3$  and stride of  $2 \times 2$  is applied to the smallest feature map to further add the scale.

Features in the modified feature maps are refined and aggregated with deformable transformers. Given a set of the modified multi-scale feature maps  $\mathbf{Z}_p = \{\mathbf{z}_i^{(p)} \mid \mathbf{z}_i^{(p)} \in \mathbb{R}^{D_p \times H'_i \times W'_i}\}_{i=1}^{L_f+1}$ , a set of refined feature maps  $\mathbf{Z}_e = \{\mathbf{z}_i^{(e)} \mid \mathbf{z}_i^{(e)} \in \mathbb{R}^{D_p \times H'_i \times W'_i}\}_{i=1}^{L_f+1}$  is obtained as  $\mathbf{Z}_e = f_{enc}(\mathbf{Z}_p, \mathbf{P})$ , where  $f_{enc}(\cdot, \cdot)$  is stacked deformable transformer encoder layers and  $\mathbf{P} = \{\mathbf{p}_i \mid \mathbf{p}_i \in \mathbb{R}^{D_p \times H'_i \times W'_i}\}_{i=1}^{L_f+1}$  is a set of multi-scale position encodings [45], which supplement the attention

modules with position and scale information to identify where each feature lies in the feature maps. The encoder helps features to acquire rich social group context by exchanging information in a feature map and between multi-scale feature maps. These enriched feature maps are fed into the deformable transformer decoder to aggregate features. Given a set of refined feature maps  $\mathbf{Z}_e$  and learnable query embeddings  $\mathbf{Q} = \{\mathbf{q}_i \mid \mathbf{q}_i \in \mathbb{R}^{2D_p}\}_{i=1}^{N_q}$ , a set of feature embeddings  $\mathbf{H} = \{\mathbf{h}_i \mid \mathbf{h}_i \in \mathbb{R}^{D_p}\}_{i=1}^{N_q}$  is obtained as  $\mathbf{H} = f_{dec}(\mathbf{Z}_e, \mathbf{Q})$ , where  $N_q$  is the number of query embeddings and  $f_{dec}(\cdot, \cdot)$  is stacked deformable transformer decoder layers. Each decoder layer predicts locations that contain features relevant to input embeddings and aggregates the features from the locations with the dynamic weighting. We design queries in such a way that one query captures at most one social group. This design enables each query to aggregate features of its target social group from the refined feature maps.

The feature embeddings are transformed into prediction results with detection heads. Here we denote the localization results in normalized image coordinates. Social group activities are recognized by predicting activities and identifying group members. The identification is performed with a group size head and group member point head. The size head predicts the number of people in a target social group, and the point head indicates group members by localizing the centers of group members' bounding boxes. This design enables our method to identify group members with simple point matching during inference as described in Sec. 3.3. The predictions of activity class probabilities  $\{\hat{\mathbf{v}}_i \mid \hat{\mathbf{v}}_i \in [0, 1]^{N_v}\}_{i=1}^{N_q}$ , group sizes  $\{\hat{s}_i \mid \hat{s}_i \in [0, 1]\}_{i=1}^{N_q}$ , and sequences of group member points  $\{\hat{\mathbf{U}}_i\}_{i=1}^{N_q}$  are obtained as  $\hat{\mathbf{v}}_i = f_v(\mathbf{h}_i)$ ,  $\hat{s}_i = f_s(\mathbf{h}_i)$ , and  $\hat{\mathbf{U}}_i = f_u(\mathbf{h}_i, \mathbf{r}_i)$ , where  $N_v$  is the number of activity classes,  $\hat{\mathbf{U}}_i = (\hat{\mathbf{u}}_j^{(i)} \mid \hat{\mathbf{u}}_j^{(i)} \in [0, 1]^2)_{j=1}^M$  is a sequence of points that indicate centers of group members' bounding boxes,  $M$  is a hyperparameter that defines the maximum group size,  $f_v(\cdot)$ ,  $f_s(\cdot)$ , and  $f_u(\cdot, \cdot)$  are the detection heads for each prediction, and  $\mathbf{r}_i \in [0, 1]^2$  is a reference point, which is used in the same way as the localization in Deformable DETR [45]. The predicted group sizes are values normalized with  $M$ . All the detection heads are composed of FFNs with subsequent sigmoid functions. We describe the details of the detection heads in the supplementary material.

Individual recognition can be performed by replacing the group recognition heads with individual recognition heads. We empirically find that using different parameters of deformable transformers for individual recognition and social group recognition does not show performance improvement and thus use shared parameters to reduce the computational costs. The details of the individual recognition heads are described in the supplementary material.

### 3.2 Loss Calculation

We view social group activity recognition as a direct set prediction problem and match predictions and ground truths with the Hungarian algorithm [22] during training following the training procedure of DETR [6]. The optimal assignment is determined by calculating the matching cost with the predicted activity class

probabilities, group sizes, and group member points. Given a ground truth set of social group activity recognition, the set is first padded with  $\phi^{(gr)}$  (no activity) to change the set size to  $N_q$ . With the padded ground truth set, the matching cost of  $i$ -th element in the ground truth set and  $j$ -th element in the prediction set is calculated as follows:

$$\mathcal{H}_{i,j}^{(gr)} = \mathbb{1}_{\{i \notin \Phi^{(gr)}\}} \left[ \eta_v \mathcal{H}_{i,j}^{(v)} + \eta_s \mathcal{H}_{i,j}^{(s)} + \eta_u \mathcal{H}_{i,j}^{(u)} \right], \quad (1)$$

$$\mathcal{H}_{i,j}^{(v)} = - \frac{\mathbf{v}_i^T \hat{\mathbf{v}}_j + (\mathbf{1} - \mathbf{v}_i)^T (\mathbf{1} - \hat{\mathbf{v}}_j)}{N_v}, \quad (2)$$

$$\mathcal{H}_{i,j}^{(s)} = |s_i - \hat{s}_j|, \quad (3)$$

$$\mathcal{H}_{i,j}^{(u)} = \frac{\sum_{k=1}^{S_i} \left\| \mathbf{u}_k^{(i)} - \hat{\mathbf{u}}_k^{(j)} \right\|_1}{S_i}, \quad (4)$$

where  $\Phi^{(gr)}$  is a set of ground-truth indices that correspond to  $\phi^{(gr)}$ ,  $\mathbf{v}_i \in \{0, 1\}^{N_v}$  is a ground truth activity label,  $s_i \in [0, 1]$  is a ground truth group size normalized with  $M$ ,  $S_i$  is an unnormalized ground truth group size,  $\mathbf{u}_k^{(i)} \in [0, 1]^2$  is a ground truth group member point normalized with the image size, and  $\eta_{\{v,s,u\}}$  are hyper-parameters. Group member points in the sequence  $\mathbf{U}_i = (\mathbf{u}_k^{(i)})_{k=1}^{S_i}$  are sorted in ascending order along  $X$  coordinates as seen from the image of the group recognition result in Fig. 2. We use this arrangement because group members are typically seen side by side at the same vertical positions in an image, and the order of group member points is clear from their positions, which makes the prediction easy. We evaluate the performances with other arrangements and compare the results in Sec. 4.4. Using Hungarian algorithm, the optimal assignment is calculated as  $\hat{\omega}^{(gr)} = \arg \min_{\omega \in \Omega_{N_q}} \sum_{i=1}^{N_q} \mathcal{H}_{i,\omega(i)}^{(gr)}$ , where  $\Omega_{N_q}$  is the set of all possible permutations of  $N_q$  elements.

The training loss for social group activity recognition  $\mathcal{L}_{gr}$  is calculated between matched ground truths and predictions as follows:

$$\mathcal{L}_v = \frac{1}{|\bar{\Phi}^{(gr)}|} \sum_{i=1}^{N_q} \left[ \mathbb{1}_{\{i \notin \Phi^{(gr)}\}} l_f(\mathbf{v}_i, \hat{\mathbf{v}}_{\hat{\omega}^{(gr)}(i)}) + \mathbb{1}_{\{i \in \Phi^{(gr)}\}} l_f(\mathbf{0}, \hat{\mathbf{v}}_{\hat{\omega}^{(gr)}(i)}) \right], \quad (5)$$

$$\mathcal{L}_s = \frac{1}{|\bar{\Phi}^{(gr)}|} \sum_{i=1}^{N_q} \mathbb{1}_{\{i \notin \Phi^{(gr)}\}} |s_i - \hat{s}_{\hat{\omega}^{(gr)}(i)}|, \quad (6)$$

$$\mathcal{L}_u = \frac{1}{|\bar{\Phi}^{(gr)}|} \sum_{i=1}^{N_q} \sum_{j=1}^{S_i} \mathbb{1}_{\{i \notin \Phi^{(gr)}\}} \left\| \mathbf{u}_j^{(i)} - \hat{\mathbf{u}}_j^{\hat{\omega}^{(gr)}(i)} \right\|_1, \quad (7)$$

where  $\lambda_{\{v,s,u\}}$  are hyper-parameters and  $l_f(\cdot, \cdot)$  is the element-wise focal loss function [28] whose hyper-parameters are described in [44].

Individual recognition is jointly learned by matching ground truths and predictions of person class probabilities, bounding boxes, and action class probabilities and calculating the losses between matched ground truths and predictions.

The matching and loss calculations are performed by slightly modifying the original matching costs and losses of Deformable DETR [45]. We describe the details of these matching and loss calculations in the supplementary material.

### 3.3 Group Member Identification

The outputs of the detection heads represent groups in group sizes and group member points that indicate centers of group members’ bounding boxes. These values have to be transformed into values that indicate individuals. We transform the predicted values into indices that refer to the elements in the individual prediction set with the following simple process during inference. To match the group member points and individual predictions, the Hungarian algorithm [22] is used instead of just calculating the closest center of a bounding box for each group member point. Hungarian algorithm can prevent multiple group member points from matching the same individuals and thus slightly improve the performance. The matching cost between  $i$ -th group member point of  $k$ -th social group prediction and  $j$ -th individual prediction is calculated as follows:

$$\mathcal{H}_{i,j}^{(gm,k)} = \frac{\|\hat{\mathbf{u}}_i^{(k)} - f_{cent}(\hat{\mathbf{b}}_j)\|_2}{\hat{c}_j}, \quad (8)$$

where  $\hat{\mathbf{b}}_j \in [0, 1]^4$  is a predicted bounding box of an individual,  $\hat{c}_j \in [0, 1]$  is a detection score of the individual, and  $f_{cent}(\cdot)$  is a function that calculates the center of a bounding box. By applying the Hungarian algorithm to this matching cost, the optimal assignment is calculated as  $\hat{\omega}^{(gm,k)} = \arg \min_{\omega \in \Omega_{N_q}} \sum_{i=1}^{\lfloor M \times \hat{s}_k \rfloor} \mathcal{H}_{i,\omega(i)}^{(gm,k)}$ , where  $\lfloor \cdot \rfloor$  rounds an input value to the nearest integer. Finally, the index set of individuals for  $k$ -th social group prediction is obtained as  $\mathbf{G}_k = \{\hat{\omega}^{(gm,k)}(i)\}_{i=1}^{\lfloor M \times \hat{s}_k \rfloor}$ .

## 4 Experiments

### 4.1 Datasets and Evaluation Metrics

We evaluate the performance of our method on two publicly available benchmark datasets: Volleyball dataset [18] and Collective Activity dataset [9]. The Volleyball dataset contains 4,830 videos of 55 volleyball matches, which are split into 3,493 training videos and 1,337 test videos. The center frame of each video is annotated with bounding boxes, actions, and one group activity. The number of action and activity classes are 9 and 8, respectively. Because the original annotations do not contain group member information, we use an extra annotation set provided by Sendo and Ukita [35]. We combine the original annotations with the group annotations in the extra set and use them for our experiments. Note that annotations other than the group annotations in the extra set are not used for a fair comparison. The Collective Activity dataset contains 44 videos of life scenes, which are split into 32 training videos and 12 test videos. The videos are annotated every ten frames with bounding boxes and actions. The

group activity is defined as the action with the largest number in the scenes. The number of action classes is 6. Because the original annotations do not have group member information, Ehsanpour *et al.* [13] annotated group labels. We use their annotations for our experiments.

We divide the evaluation into two parts: group activity recognition and social group activity recognition. In the evaluation of group activity recognition, we follow the detection-based settings [5, 13, 33, 41] and use classification accuracy as an evaluation metric. Because our method is designed to predict multiple group activities, we need to select one from them for group activity recognition. We choose the predicted activity of the highest probability and compare it with the ground truth activity. In the evaluation of social group activity recognition, different metrics are used for each dataset because each scene in the Volleyball dataset contains only one social group activity, while that in the Collective Activity dataset contains multiple social group activities. For the Volleyball dataset, group identification accuracy is used as an evaluation metric. One group prediction is first selected in the same way as group activity recognition, and then the predicted bounding boxes of the group members are compared with the ground truth boxes. The selected prediction results are correct if the predicted activity is correct and the predicted boxes have IoUs larger than 0.5 with the corresponding ground truth boxes. For the Collective Activity dataset, mAP is used as an evaluation metric. Prediction results are judged as true positives if the predicted activities are correct, and all the predicted boxes of the group members have IoUs larger than 0.5 with the corresponding ground truth boxes.

## 4.2 Implementation Details

We use the RGB stream of I3D [7] as a backbone feature extractor and input features from *Mixed\_3c*, *Mixed\_4f*, and *Mixed\_5c* layers into the deformable transformers. The hyper-parameters of the deformable transformers are set in accordance with the setting of Deformable DETR [45], where  $L_f = 3$ ,  $D_p = 256$ , and  $N_q = 300$ . We initialize I3D with the parameters trained on the Kinetics dataset [19] and deformable transformers with the parameters trained on the COCO dataset [29]. We use the AdamW [30] optimizer with the batch size of 16, the initial learning rate of  $10^{-4}$ , and the weight decay of  $10^{-4}$ . Training epochs are set to 120, and the learning rate is decayed after 100 epochs. We set the length of the sequence  $T$  to 9. Ground truth labels of the center frame are used to calculate the losses. To augment the training data, we randomly shift frames in the temporal direction and use bounding boxes from visual trackers as ground truth boxes when a non-annotated frame is at the center. We also augment the training data by random horizontal flipping, scaling, and cropping. Following the DETR’s training [6], auxiliary losses are used to boost the performance. The maximum group size  $M$  is set to 12. The hyper-parameters are set as  $\eta_v = \lambda_v = 2$ ,  $\eta_s = \lambda_s = 1$ , and  $\eta_u = \lambda_u = 5$ .

While evaluating performances with the Collective Activity dataset, some specific settings are used. For the evaluation of group activity recognition, training epochs are set to 10, and the learning rate is decayed after 5 epochs because

Table 1: Comparison against state-of-the-art methods on group activity recognition. The values with and without the brackets demonstrate the performances in the ground-truth-based and detection-based settings, respectively. The performances of individual action recognition are shown for future reference.

Method	Volleyball		Collective Activity			
	Activity	Action	Activity	Action	Activity	Action
SSU [5]	86.2 (90.6)	– (81.8)	– (–)	– (–)	– (–)	– (–)
stagNet [33]	87.6 (89.3)	– (–)	87.9 (89.1)	– (–)	– (–)	– (–)
ARG [41]	91.5 (92.5)	39.8 (83.0)	86.1 (88.1)	49.6 (77.3)	– (–)	– (–)
CRM [4]	– (93.0)	– (–)	– (85.8)	– (–)	– (–)	– (–)
PRL [16]	– (91.4)	– (–)	– (–)	– (–)	– (–)	– (–)
Actor-Transformers [14]	– (94.4)	– (85.9)	– (92.8)	– (–)	– (–)	– (–)
Ehsanpour <i>et al.</i> [13]	93.0 (93.1)	41.8 (83.3)	89.4 (89.4)	55.9 (78.3)	– (–)	– (–)
Pramono <i>et al.</i> [32]	– (95.0)	– (83.1)	– (95.2)	– (–)	– (–)	– (–)
DIN [42]	– (93.6)	– (–)	– (95.9)	– (–)	– (–)	– (–)
GroupFormer [26]	95.0* (95.7)	– (85.6)	85.2* (87.5 <sup>†</sup> /96.3)	– (–)	– (–)	– (–)
Ours	<b>96.0</b> (–)	<b>65.0</b> (–)	<b>96.5</b> (–)	<b>64.9</b> (–)	– (–)	– (–)

\* We evaluated the performance with the publicly available source codes.

† We evaluated but were not able to reproduce the reported accuracy because the configuration file for the Collective Activity dataset is not publicly available.

the losses converge in a few epochs due to the limited diversity of the scenes in the dataset. For the evaluation of social group activity recognition, the length of the sequence  $T$  is set to 17 following the setting of Ehsanpour *et al.* [13].

### 4.3 Group Activity Recognition

**Comparison against State-of-the-Art.** We compare our method against state-of-the-art methods on group activity recognition. Table 1 shows the comparison results. The values without the brackets demonstrate the detection-based performances, while those inside the brackets indicate the performances with ground truth bounding boxes. We show the performances of individual action recognition for future reference. Several detection-based performances are not reported because existing works typically use ground-truth boxes for the evaluation. To compare the effectiveness with these methods, we evaluate GroupFormer [26], which is the strongest baseline of group activity recognition, with predicted boxes of Deformable DETR [45]. Note that Deformable DETR is fine-tuned on each dataset for a fair comparison, which demonstrates 90.8 and 90.2 mAP on the Volleyball and Collective Activity datasets, respectively.

As seen from the table, our method outperforms state-of-the-art methods in the detection-based setting. We confirm that GroupFormer shows the performance degradation as well as the previous methods [5, 13, 33, 41] when predicted bounding boxes are used. These results indicate that the latest region-feature-based method still suffers from incomplete person localization and that our

Table 2: Analysis on the effect of the group annotations with the Volleyball dataset. The values with and without the brackets demonstrate the performances in the ground-truth-based and detection-based settings, respectively.

Method	Annotation type	Activity
GroupFormer [26]	Original	95.0* (95.7)
	Group	93.2 <sup>‡</sup> (96.1*)
Ours	Original	95.0 ( - )
	Group	<b>96.0</b> ( - )

\* We evaluated the performance with the publicly available source codes.

<sup>‡</sup> We trained a group member detector and evaluated the performance with publicly available source codes.

feature generation has advantages over these methods. Even compared to the ground-truth-based performances, our method shows the best performance. It is worth noting that our method uses only RGB images as inputs, while GroupFormer utilizes optical flows and pose information in addition to RGB data. These results suggest that features generated by our method are more effective than region features and that it is not optimal to restrict regions of features to bounding boxes.

**Analysis on Group Annotations.** As described in Sec. 4.1, we use the additional group annotations to fully leverage our social group activity recognition capability. We analyze the effect of the group annotations on group activity recognition by investigating the performances of both GroupFormer [26] and our method with and without the group annotations. Note that hereinafter we use the Volleyball dataset for analyses because the diversity of the scenes in the Collective Activity dataset is limited. To evaluate GroupFormer with the group annotations in the detection-based setting, we trained Deformable DETR [45] with bounding boxes of only group members, which is intended to detect only people involved in activities. The detector shows the performance of 87.1 mAP. Among all the results, GroupFormer with the group annotations in the ground-truth-based setting demonstrates the best performance. However, the performance is substantially degraded when the predicted boxes are used. This is probably because group member detection underperforms and degrades the recognition performance. As our method does not rely on bounding boxes to predict group activities, the performance does not degrade even if group members cannot be identified correctly. Accordingly, our method demonstrates the best performance in the detection-based setting.

#### 4.4 Social Group Activity Recognition

**Comparison against State-of-the-Art.** To demonstrate the effectiveness of our method on social group activity recognition, we compare our method against

Table 3: Comparison against state-of-the-art social group activity recognition methods with the Volleyball dataset.

Method	Accuracy	Right				Left			
		Set	Spike	Pass	Winpoint	Set	Spike	Pass	Winpoint
Ehsanpour <i>et al.</i> [13] <sup>§</sup>	44.5	17.2	<b>74.0</b>	49.0	29.9	19.7	<b>79.6</b>	25.0	28.4
GroupFormer [26] <sup>‡</sup>	48.8	25.0	56.6	59.0	<b>51.7</b>	31.5	55.3	58.8	51.0
Ours	<b>60.6</b>	<b>35.9</b>	68.2	<b>81.9</b>	50.6	<b>50.6</b>	53.6	<b>74.3</b>	<b>56.9</b>

<sup>§</sup> Because the source codes are not publicly available, we implemented their algorithm based on our best understanding and evaluated the performance.

<sup>‡</sup> We trained a group member detector and evaluated the performance with publicly available source codes.

Table 4: Comparison against a state-of-the-art social group activity recognition method with the Collective Activity dataset.

Method	mAP	Crossing	Waiting	Queueing	Walking	Talking
Ehsanpour <i>et al.</i> [13]	<b>51.3</b>	–	–	–	–	–
Ours	46.0	49.2	64.5	54.1	55.6	6.56

Ehsanpour *et al.*’s method [13], which is a state-of-the-art method that tackles social group activity recognition, and GroupFormer [26], which is the strongest baseline on group activity recognition. Due to the unavailability of both Ehsanpour *et al.*’s source codes and their performance report on the Volleyball dataset, we implemented their algorithm based on our best understanding and evaluated the performance on the dataset. For the evaluation of GroupFormer, we trained Deformable DETR [45] in the same way as described in the group annotation analysis section for detecting group members. Because this group member detection cannot be applied to multiple social groups, we evaluate GroupFormer only on the Volleyball dataset.

Table 3 shows the results on the Volleyball dataset. As shown in the table, our method yields significant performance gains over the other methods, which demonstrates the improvement on group member identification as well as on activity recognition. Our method aggregates features that are embedded with clues for grouping people from feature maps. It is highly likely that this feature aggregation contributes to the high accuracy of identifying activities with different distributions of group members in an image. We qualitatively analyze how features are aggregated depending on the distribution of group members and discuss the analysis results towards the end in the qualitative analysis section.

The comparison results on the Collective Activity dataset are listed in Table 4. As seen from the table, Ehsanpour *et al.*’s method shows better performance than our method. We find that our method demonstrates relatively low performance on the activity “Talking”. This low performance is probably attributed to the number of samples in training data. In the test data, 86% of

Table 5: Analysis on group sizes with Volleyball dataset.

Method	Group size (Training data ratio)					
	1 (36%)	2 (21%)	3 (19%)	4 (6%)	5 (5%)	6 (12%)
Ehsanpour <i>et al.</i> [13] <sup>§</sup>	45.3	<b>48.2</b>	<b>61.2</b>	27.3	15.8	32.5
GroupFormer [26] <sup>‡</sup>	57.3	29.6	58.4	<b>28.4</b>	<b>44.7</b>	54.4
Ours	<b>83.6</b>	42.9	52.4	26.1	39.5	<b>63.8</b>

<sup>§</sup> Because the source codes are not publicly available, we implemented their algorithm based on our best understanding and evaluated the performance.

<sup>‡</sup> We trained a group member detector and evaluated the performance with publicly available source codes.

Table 6: Analysis on the order of member points with the Volleyball dataset.

Order of the group member points	Probability of changes in order	Accuracy
Ascending order in X coordinates	7.4%	<b>60.6</b>
Ascending order in Y coordinates	13%	55.5

samples with the activity ‘‘Talking’’ have the group sizes of four, while the training data has only 57 samples whose group sizes are four, which is 0.8% of the training data. As our method learns to predict group sizes, the number of samples in training data for each group size affects the performance. We analyze this effect in the subsequent section.

**Analysis on Group Sizes.** The group size prediction is one of the key factors to identify group members and thus affects social group activity recognition performance. To analyze this effect, we evaluate the performance on each group size and compare the results with Ehsanpour *et al.*’s method [13] and GroupFormer [26]. Table 5 shows the results. As shown in the table, the performances of our method are moderately correlated to the training data ratios, while the other two methods do not show the correlation. This is the drawback of our method that relies on group size learning. However, our method shows the competitive performances on both small and large group sizes if there are a certain amount of training data. In contrast, each of the other two methods shows the competitive performances only on either large or small group sizes. These results imply that our method does not have the performance dependence on group sizes and thus can achieve high performance with large-scale training data.

**Analysis on Order of Group Member Points.** As described in Sec. 3.2, group member points in a ground truth point sequence are sorted in ascending order along  $X$  coordinates. To confirm the effectiveness of this arrangement, we compare the performances with two arrangements. Table 6 shows the comparison results. As shown in the table, our method demonstrates better performance when group member points are sorted in ascending order along  $X$  coordinates

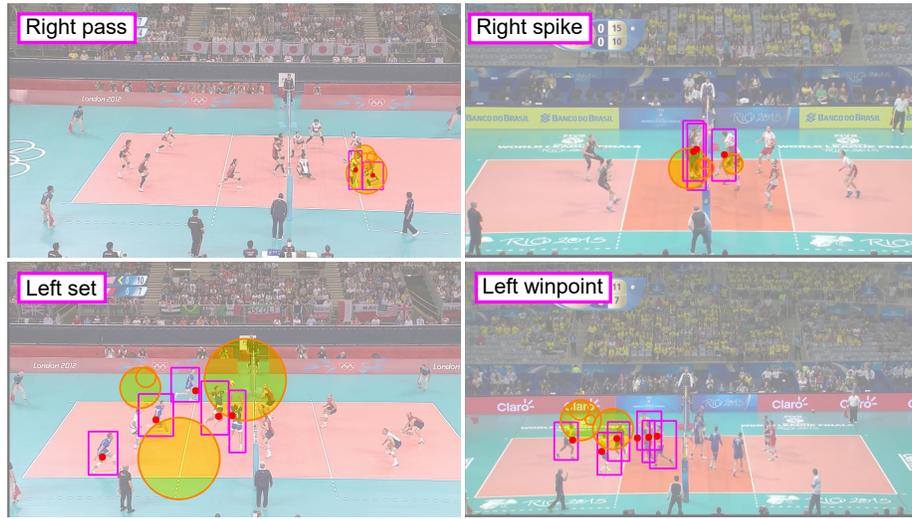


Fig. 3: Visualization of the attention locations in the deformable transformer decoder. We show the locations of the top four attention weights. The large circles mean that the locations are in the low-resolution feature maps.

than in ascending order along  $Y$  coordinates. The probabilities in the table indicate the ratio of the changes of the point order when small perturbations are added to ground-truth bounding box positions. The higher probability implies that the order of group member points changes more frequently when group members move. These results suggest that the order changes more frequently when group member points are sorted in ascending order along  $Y$  coordinates and that the order is difficult to predict with slight differences of box positions.

**Qualitative Analysis.** The deformable attention modules are the critical components to aggregate features relevant to social group activity recognition and generate social group features. To analyze how the attention modules aggregate features for various social group activities, we visualize the attention locations of the transformer decoder in Fig. 3. We show locations with the top four attention weights in the last layer of the decoder. The purple bounding boxes show the group members, the red circles show the predicted group member points, and the yellow circles show the attention locations. The small and large yellow circles mean that the locations are in the high and low-resolution feature maps, respectively, showing a rough range of image areas affecting the generated features. The figure shows that features are typically aggregated from low-resolution feature maps if group members are located in broad areas, and vice versa. These results indicate that the attention modules can effectively aggregate features depending on the distribution of group members and contribute to improving the performance of social group activity recognition.

## 5 Conclusions

We propose a novel social group activity recognition method that leverages deformable transformers to generate effective social group features. This feature generation obviates the need for region features and hence makes the effectiveness of the social group features person-localization-agnostic. Furthermore, the group member information extracted from the features is represented so concisely that our method can identify group members with simple Hungarian matching, resulting in high-performance social group activity recognition. We perform extensive experiments and show significant improvement over existing methods.

## Acknowledgement

Computational resource of AI Bridging Cloud Infrastructure (ABCI) provided by National Institute of Advanced Industrial Science and Technology (AIST) was used.

## References

1. Amer, M.R., Lei, P., Todorovic, S.: HiRF: Hierarchical random field for collective activity recognition in videos. In: ECCV (September 2014)
2. Amer, M.R., Todorovic, S.: Sum product networks for activity recognition. IEEE TPAMI **38**(4), 800–813 (April 2016)
3. Amer, M.R., Todorovic, S., Fern, A., Zhu, S.C.: Monte carlo tree search for scheduling activity recognition. In: ICCV (December 2013)
4. Azar, S.M., Atigh, M.G., Nickabadi, A., Alahi, A.: Convolutional relational machine for group activity recognition. In: CVPR (June 2019)
5. Bagautdinov, T.M., Alahi, A., Fleuret, F., Fua, P.V., Savarese, S.: Social scene understanding: End-to-end multi-person action localization and collective activity recognition. In: CVPR (July 2017)
6. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (November 2020)
7. Carreira, J., Zisserman, A.: Quo vadis, action recognition? A new model and the kinetics dataset. In: CVPR (July 2017)
8. Cho, K., van Merriënboer, B., Çaglar Gülçehre, Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. In: EMNLP (October 2014)
9. Choi, W., Shahid, K., Savarese, S.: What are they doing? : Collective activity classification using spatio-temporal relationship among people. In: ICCVW (September 2009)
10. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: ICCV (October 2017)
11. Dai, X., Chen, Y., Yang, J., Zhang, P., Yuan, L., Zhang, L.: Dynamic DETR: End-to-end object detection with dynamic attention. In: ICCV (October 2021)
12. Deng, Z., Vahdat, A., Hu, H., Mori, G.: Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition. In: CVPR (June 2016)

13. Ehsanpour, M., Abedin, A., Saleh, F., Shi, J., Reid, I., Rezatofghi, H.: Joint learning of social groups, individuals action and sub-group activities in videos. In: ECCV (August 2020)
14. Gavriluyk, K., Sanford, R., Javan, M., Snoek, C.G.M.: Actor-transformers for group activity recognition. In: CVPR (June 2020)
15. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* **9**(8), 1735–1780 (November 1997)
16. Hu, G., Cui, B., He, Y., Yu, S.: Progressive relation learning for group activity recognition. In: CVPR (June 2020)
17. Ibrahim, M.S., Mori, G.: Hierarchical relational networks for group activity recognition and retrieval. In: ECCV (September 2018)
18. Ibrahim, M.S., Muralidharan, S., Deng, Z., Vahdat, A., Mori, G.: A hierarchical deep temporal model for group activity recognition. In: CVPR (June 2016)
19. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., Suleyman, M., Zisserman, A.: The kinetics human action video dataset (May 2017), arXiv:1705.06950
20. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: ICLR (April 2017)
21. Kong, L., Qin, J., Huang, D., Wang, Y., Gool, L.V.: Hierarchical attention and context modeling for group activity recognition. In: ICASSP (April 2018)
22. Kuhn, H.W., Yaw, B.: The hungarian method for the assignment problem. *Naval Res. Logist. Quart* pp. 83–97 (March 1955)
23. Lan, T., Sigal, L., Mori, G.: Social roles in hierarchical models for human activity recognition. In: CVPR (June 2012)
24. Lan, T., Wang, Y., Yang, W., Mori, G.: Beyond actions: Discriminative models for contextual group activities. In: NIPS (December 2010)
25. Lan, T., Wang, Y., Yang, W., Robinovitch, S.N., Mori, G.: Discriminative latent models for recognizing contextual group activities. *IEEE TPAMI* **34**(8), 1549–1562 (August 2012)
26. Li, S., Cao, Q., Liu, L., Yang, K., Liu, S., Hou, J., Yi, S.: GroupFormer: Group activity recognition with clustered spatial-temporal transformer. In: ICCV (October 2021)
27. Li, X., Chuah, M.C.: SBGAR: Semantics based group activity recognition. In: ICCV (October 2017)
28. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (October 2017)
29. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: ECCV (September 2014)
30. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (May 2019)
31. Ng, A., Jordan, M., Weiss, Y.: On spectral clustering: Analysis and an algorithm. In: NIPS (December 2002)
32. Pramono, R.R.A., Chen, Y.T., Fang, W.H.: Empowering relational network by self-attention augmented conditional random fields for group activity recognition. In: ECCV (August 2020)
33. Qi, M., Qin, J., Li, A., Wang, Y., Luo, J., Gool, L.V.: stagNet: An attentive semantic rnn for group activity recognition. In: ECCV (September 2018)
34. Rezatofghi, H., Tsoi, N., Gwak, J., Sadeghian, A., Reid, I., Savarese, S.: Generalized intersection over union: A metric and a loss for bounding box regression. In: CVPR (June 2019)

35. Sendo, K., Ukita, N.: Heatmapping of people involved in group activities. In: MVA (May 2019)
36. Shu, T., Todorovic, S., Zhu, S.C.: CERN: Confidence-energy recurrent network for group activity recognition. In: CVPR (July 2017)
37. Sun, Z., Cao, S., Yang, Y., Kitani, K.M.: Rethinking transformer-based set prediction for object detection. In: ICCV (October 2021)
38. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NIPS (December 2017)
39. Veličkovič, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: ICLR (April 2018)
40. Wang, M., Ni, B., Yang, X.: Recurrent modeling of interaction context for collective activity recognition. In: CVPR (July 2017)
41. Wu, J., Wang, L., Wang, L., Guo, J., Wu, G.: Learning actor relation graphs for group activity recognition. In: CVPR (June 2019)
42. Yuan, H., Ni, D., Wang, M.: Spatio-temporal dynamic inference network for group activity recognition. In: ICCV (October 2021)
43. Zhou, H., Kadav, A., Shamsian, A., Geng, S., Lai, F., Zhao, L., Liu, T., Kapadia, M., Graf, H.P.: COMPOSER: Compositional learning of group activity in videos (December 2021), arXiv:2112.05892
44. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points (April 2019), arXiv:1904.07850
45. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: Deformable transformers for end-to-end object detection. In: ICLR (May 2021)

## Supplementary Material

### A Individual Recognition

In our method, individuals are recognized by simply adding an action classification head to the detection heads in Deformable DETR [45]. Given a set of feature embedding  $\mathbf{H} = \{\mathbf{h}_i \mid \mathbf{h}_i \in \mathbb{R}^{D_p}\}_{i=1}^{N_q}$  from the deformable transformer decoder, the predictions of person class probabilities  $\{\hat{c}_i \mid \hat{c}_i \in [0, 1]\}_{i=1}^{N_q}$ , bounding boxes  $\{\hat{\mathbf{b}}_i \mid \hat{\mathbf{b}}_i \in [0, 1]^4\}_{i=1}^{N_q}$ , and action class probabilities  $\{\hat{\mathbf{a}}_i \mid \hat{\mathbf{a}}_i \in [0, 1]^{N_a}\}_{i=1}^{N_q}$  are obtained as  $\hat{c}_i = f_c(\mathbf{h}_i)$ ,  $\hat{\mathbf{b}}_i = f_b(\mathbf{h}_i, \mathbf{r}_i)$ , and  $\hat{\mathbf{a}}_i = f_a(\mathbf{h}_i)$ , where  $N_q$  is the number of query embeddings,  $N_a$  is the number of action classes,  $f_c(\cdot)$ ,  $f_b(\cdot, \cdot)$ , and  $f_a(\cdot)$  are the detection heads for the predictions, and  $\mathbf{r}_i \in [0, 1]^2$  is a reference point, which is used in the same way as the localization in Deformable DETR. Note that the localization results are denoted in the normalized image coordinates.

We view individual recognition as a direct set prediction problem and match predictions and ground truths with the Hungarian algorithm [22] during training. The optimal assignment of ground truths and predictions is determined by calculating the matching cost with the predicted person class probabilities, bounding boxes, and action class probabilities. Given a ground truth set of individual recognition, the set is first padded with  $\phi^{(id)}$  (no person) to change the size of the set to  $N_q$ . Using the padded ground truth set, the matching cost of  $i$ -th element in the ground truth set and  $j$ -th element in the prediction set for individual recognition is calculated as follows:

$$\mathcal{H}_{i,j}^{(id)} = \mathbb{1}_{\{i \notin \Phi^{(id)}\}} \left[ \eta_c \mathcal{H}_{i,j}^{(c)} + \eta_b \mathcal{H}_{i,j}^{(b)} + \eta_o \mathcal{H}_{i,j}^{(o)} + \eta_a \mathcal{H}_{i,j}^{(a)} \right], \quad (9)$$

$$\mathcal{H}_{i,j}^{(c)} = -\hat{c}_j, \quad (10)$$

$$\mathcal{H}_{i,j}^{(b)} = \left\| \mathbf{b}_i - \hat{\mathbf{b}}_j \right\|_1, \quad (11)$$

$$\mathcal{H}_{i,j}^{(o)} = -f_{GIoU}(\mathbf{b}_i, \hat{\mathbf{b}}_j), \quad (12)$$

$$\mathcal{H}_{i,j}^{(a)} = -\left( \frac{\mathbf{a}_i^T \hat{\mathbf{a}}_j + (\mathbf{1} - \mathbf{a}_i)^T (\mathbf{1} - \hat{\mathbf{a}}_j)}{N_a} \right), \quad (13)$$

where  $\Phi^{(id)}$  is a set of ground-truth indices that correspond to  $\phi^{(id)}$ ,  $\mathbf{b}_i \in [0, 1]^4$  is a ground truth bounding box normalized with the image size,  $\mathbf{a}_i \in \{0, 1\}^{N_a}$  is a ground truth action label,  $f_{GIoU}(\cdot, \cdot)$  is a function that calculates generalized IoU [34], and  $\eta_{\{c,b,o,a\}}$  are the hyper-parameters. The Hungarian algorithm is applied to the matching cost to find the optimal assignment  $\hat{\omega}^{(id)} = \arg \min_{\omega \in \Omega_{N_q}} \sum_{i=1}^{N_q} \mathcal{H}_{i,\omega(i)}^{(id)}$ , where  $\Omega_{N_q}$  is the set of all possible permutations of  $N_q$  elements.

The training loss for individual recognition  $\mathcal{L}_{id}$  is calculated between matched ground truths and predictions as follows:

$$\mathcal{L}_{id} = \lambda_c \mathcal{L}_c + \lambda_b \mathcal{L}_b + \lambda_o \mathcal{L}_o + \lambda_a \mathcal{L}_a, \quad (14)$$

$$\mathcal{L}_c = \frac{1}{|\Phi^{(id)}|} \sum_{i=1}^{N_q} [\mathbb{1}_{\{i \notin \Phi^{(id)}\}} l_f([1], [\hat{c}_{\omega^{(id)}(i)}]) + \mathbb{1}_{\{i \in \Phi^{(id)}\}} l_f([0], [\hat{c}_{\omega^{(id)}(i)}])], \quad (15)$$

$$\mathcal{L}_b = \frac{1}{|\Phi^{(id)}|} \sum_{i=1}^{N_q} \mathbb{1}_{\{i \notin \Phi^{(id)}\}} \left\| \mathbf{b}_i - \hat{\mathbf{b}}_{\omega^{(id)}(i)} \right\|_1, \quad (16)$$

$$\mathcal{L}_o = \frac{1}{|\Phi^{(id)}|} \sum_{i=1}^{N_q} \mathbb{1}_{\{i \notin \Phi^{(id)}\}} \left[ 1 - f_{GIoU}(\mathbf{b}_i, \hat{\mathbf{b}}_{\omega^{(id)}(i)}) \right], \quad (17)$$

$$\mathcal{L}_a = \frac{1}{|\Phi^{(id)}|} \sum_{i=1}^{N_q} \mathbb{1}_{\{i \notin \Phi^{(id)}\}} l_f(\mathbf{a}_i, \hat{\mathbf{a}}_{\omega^{(id)}(i)}), \quad (18)$$

where  $\lambda_{\{c,b,o,a\}}$  are hyper-parameters and  $l_f(\cdot, \cdot)$  is the element-wise focal loss function [28] whose hyper-parameters are described in [44].

In our training, the hyper-parameters  $\eta_{\{c,b,o,a\}}$  and  $\lambda_{\{c,b,o,a\}}$  are set as  $\eta_c = \lambda_c = 1$ ,  $\eta_b = \lambda_b = 5$ ,  $\eta_o = \lambda_o = 2$ , and  $\eta_a = \lambda_a = 2$ .

## B Implementation Details of Detection Heads

In our method, all the detection heads are constituted by feed-forward networks with the subsequent sigmoid functions. The details of the detection heads are as follows:

### Person class head

This head has 1 linear layer with the subsequent sigmoid function.

### Box head

This head has 3 linear layers with the ReLU activation between the layers and the subsequent sigmoid function. A reference point is added to each corresponding box position before applying the sigmoid function.

### Action head

This head has 1 linear layer with the subsequent sigmoid function.

### Activity head

This head has 1 linear layer with the subsequent sigmoid function.

### Group size head

This head has 3 linear layers with the ReLU activation between the layers and the subsequent sigmoid function.

### Member point head

This head has 3 linear layers with the ReLU activation between the layers and the subsequent sigmoid function.  $2 \times M$  values are output from the last linear layer and then split into  $M$  group member points, where  $M$  denotes the maximum group size. A reference point is added to each corresponding group member point before applying the sigmoid function.

## C Group Annotations in The Volleyball Dataset

Group annotations are critical components to fully leverage the learning capability of our method. In the evaluation of the Volleyball dataset [18], we use the original annotation set combined with the extra annotation set provided by Sendo and Ukita [35] because the original annotations do not contain group information. The group annotations in the extra set are transferred to the original set by matching bounding boxes from each set with intersection over union (IoU). IoU is first calculated for each pair of a box from the original set and that from the extra set in the same frame. The calculated IoU values are then used as costs for the Hungarian algorithm [22] to match the boxes. If a box from the extra set has a label indicating that the person in the box is involved in an activity, we assign a group member flag to the matched box from the original set.

The players involved in each group activity are defined by Sendo and Ukita [35] as follows:

### Pass

Players who are trying an underhand pass independently of whether or not they successfully do it.

### Set

A player who is doing an overhand pass and those who will spike the ball whether they are trying or faking.

### Spike

Players who are spiking and blocking.

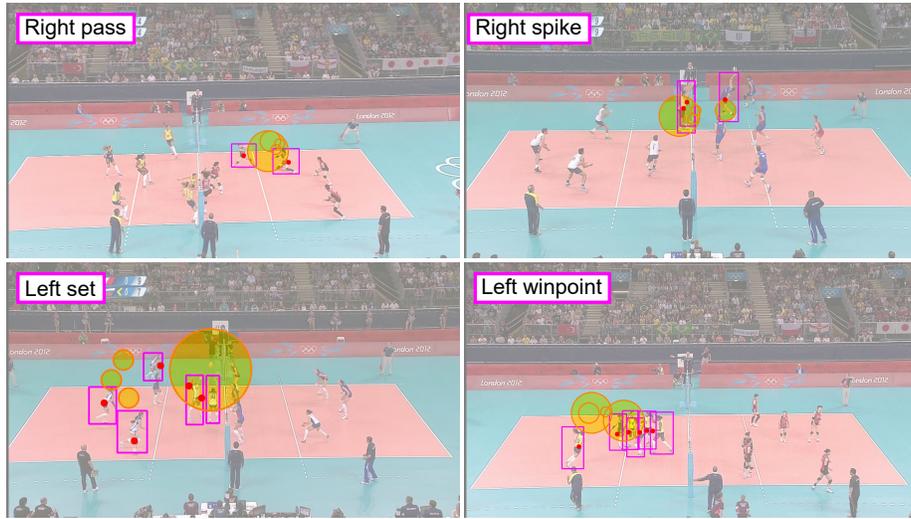
### Winpoint

All players in the team scoring a point. This group activity is observed for a few seconds right after scoring.

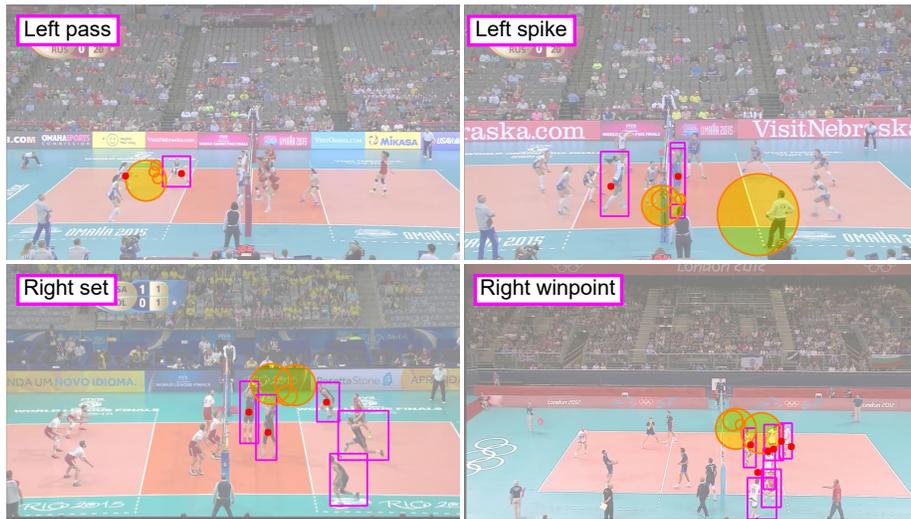
## D Additional Qualitative Analysis

We further analyze the recognition results qualitatively with our method’s success and failure cases on the Volleyball dataset [18]. The results of the successful cases and failure cases are shown in Fig. 4a and 4b, respectively. The purple bounding boxes show the ground truth group members, the red circles show the predicted group member points, and the yellow circles show the attention locations. The small and large yellow circles mean that the locations are in the high and low-resolution feature maps, respectively, offering a rough range of image areas affecting the features used for the predictions.

As seen from the figures, features are successfully aggregated from the areas around the group members in the successful cases, while those are aggregated from the regions around the non-group members, backgrounds, and part of the group members in the failure cases. It is worth noting that our method successfully recognizes social group activities even when one group member is apart from the other members such as the cases of “Right spike” and “Left winpoint” in Fig. 4a, demonstrating the effectiveness of our feature aggregation method.



(a) Successful cases.



(b) Failure cases.

Fig. 4: Visualization of the social group activity recognition results. The purple bounding boxes, red circles, and yellow circles show the ground truth group members, predicted group member points, and attention locations in the deformable transformer decoder, respectively.

In the failure case of “Left pass”, a non-group member is falsely recognized as a group member probably because the non-group member has the pose of the underhand pass, which is quite similar to the group member. In the failure case of “Left spike”, a group member cannot be identified due to the occlusion. To correctly identify group members in these cases, long-term temporal context should be leveraged effectively. In the failure cases of “Right set” and “Right winpoint”, the group members are widely distributed especially in the vertical direction. As discussed in the main manuscript, the group member point prediction is designed on the assumption that group members are seen side by side at the same vertical positions in an image. This design might affect the performance of the failure cases. These observations present opportunities for future work.