# Learning Semantics-Preserving Attention and Contextual Interaction for Group Activity Recognition

Yansong Tang<sup>®</sup>, Student Member, IEEE, Jiwen Lu<sup>®</sup>, Senior Member, IEEE, Zian Wang<sup>®</sup>, Ming Yang, Member, IEEE, and Jie Zhou, Senior Member, IEEE

Abstract—In this paper, we investigate the problem of group activity recognition by learning semantics-preserving attention 2 and contextual interaction among different people. Conventional 3 methods usually aggregate the features extracted from individual persons by pooling operations, which lack physical meaning 5 and cannot fully explore the contextual information for group 6 activity recognition. To address this, we develop a Semantics-Preserving Teacher-Student (SPTS) networks architecture. Our SPTS networks first learn a Teacher Network in the semantic 9 domain that classifies the word of group activity based on the 10 words of individual actions. Then, we design a Student Network 11 in the appearance domain that recognizes the group activity 12 according to the input video. We enforce the Student Network 13 to mimic the Teacher Network in the learning procedure. In this 14 way, we allocate semantics-preserving attention to different 15 people, which is more effective to seek the key people and 16 discard the misleading people, while no extra labeled data are 17 required. Moreover, a group of people inherently lie in a graph-18 based structure, where the people and their relationship can 19 be regarded as the nodes and edges of a graph, respectively. 20 Based on this, we build two graph convolutional modules on 21 both the Teacher Network and the Student Network to reason the 22 dependency among different people. Furthermore, we extend our 23 approach on action segmentation task based on its intermediate 24 features. The experimental results on four datasets for group 25 activity analysis clearly show the superior performance of our 26 method in comparison with the state-of-the-art. 27

Index Terms—Semantics-preserving, attention, group activity
 recognition, Teacher-Student networks.

30

31

32

# I. INTRODUCTION

**G**ROUP activity recognition (*a.k.a.* collective activity recognition), which refers to discerning what a group

Manuscript received September 3, 2018; revised March 2, 2019 and April 26, 2019; accepted April 29, 2019. This work was supported in part by the National Natural Science Foundation of China under Grant U1813218, Grant 61822603, Grant U1713214, Grant 61672306, and Grant 61572271. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Tolga Tasdizen. (*Corresponding author: Jiwen Lu.*)

Y. Tang, J. Lu, Z. Wang, and J. Zhou are with the State Key Laboratory of Intelligent Technologies and Systems, Beijing Research Center for Information Science and Technology (BNRist), Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: tys15@mails.tsinghua.edu.cn; lujiwen@tsinghua.edu.cn; wza15@mails.tsinghua.edu.cn; jzhou@ tsinghua.edu.cn).

M. Yang is with Horizon Robotics, Inc., Beijing 100080, China (e-mail: ming.yang@horizon-robotics.com).

This paper has supplementary downloadable material available at http://ieeexplore.ieee.org, provided by the authors.

Digital Object Identifier 10.1109/TIP.2019.2914577

of people are doing in a video, has attracted growing attention in the realm of computer vision over the past decade [1]–[7]. There are wide real-world applications for group activity recognition including traffic surveillance, social role understanding and sports video analysis. Compared with conventional action recognition which focuses on a single person, group activity recognition is a more challenging task as it requires further understanding of high-level relationships among different people. Hence, it is desirable to design a model to aggregate the individual dynamics across people and exploit their contextual information for effective group activity recognition.

Over the past few years, great efforts have been devoted to 45 mining the contextual information for group activity recog-46 nition. In the early period, a typical series of approaches 47 are developed to design graph-based structure models based 48 on hand-crafted features [7]-[10]. However, these methods 49 require strong prior knowledge and lack discriminative power 50 to model the temporal evolution of group activity. In recent 51 years, with the spectacular progress of deep learning methods, 52 researchers have attempted to build different deep neural 53 networks [2], [3] for group activity recognition. Most of these 54 methods treat all participants with equal importance, and 55 integrate the features of individual actions by simple pooling 56 operators. However, the group activity is usually sensitive to a 57 few key persons, whose actions essentially define the activity, 58 and other people may bring ambiguous information and mis-59 lead the recognition process. Let's take Fig. 1 as an example. 60 The bottom of Fig. 1 shows a frame sampled from a video 61 clip in Volleyball dataset [2]. Obviously, the "spiking" person 62 shall provide more discriminative information for recognizing 63 the "right spike" activity, and those "standing" people may 64 bring some confounding information. To address these, several 65 attention-based methods [5], [11] have been proposed to assign 66 different weights to different people. Specifically, the weights 67 are learned based on the features extracted from input videos, 68 and are allocated to their corresponding features. However, 69 such a "self-attention" scheme essentially lacks physical expla-70 nation and is not reliable enough to find the key person for 71 activity recognition. 72

In this work, we move a new step towards the interaction <sup>73</sup> of appearance domain and semantic domain, and propose <sup>74</sup> a Semantics-Preserving Teacher-Student (SPTS) model for <sup>75</sup>

1057-7149 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

33

34

35

36

37

38

39

40

41

42

43



Fig. 1. The basic idea of the SPTS networks. In the semantic domain, the task is to map the *words* of individual actions, which can be treated as a caption of the video [4], to the word of group activity. In the appearance domain, we attempt to predict the label of group activity based on the corresponding input video. We first learn a Teacher Network in the semantic domain, and then employ the learned attention information, which represents the different importance of different people for recognizing the group activity, to guide a Student Network in the appearance domain. (Best viewed in color.)

group activity recognition. Fig. 1 shows the basic idea of 76 our approach. Concretely, we first learn a high-performance 77 model with typical attention mechanism (namely Teacher 78 Network) to map the individual actions to group activity in 79 the semantic domain. Next, we develop another model (namely 80 Student Network), which predicts the group activity from the 81 individual actions in the appearance domain. Then, we design 82 a unified framework to utilize the attention knowledge in the 83 Teacher Network to guide the Student Network. As the inputs 84 of our Teacher Network are generated from the off-the-shelf 85 single-action labels, our method requires no extra labelled 86 data and only takes additional 2.70% computational time cost. 87 Moreover, most conventional methods model the features of 88 group people as regular tensor-based vectors, which ignore 89 the intrinsic dependency among different people. To address 90 this, we construct two types of graphs in semantic domain 91 and appearance domain, respectively. The nodes of the graph 92 contain the extracted features of the individual persons, while 93 the adjacency matrices that encode their spatial coordinates are 94 used to describe the relationship among different people. Since 95 the graph of features lies in a non-Euclidean space, we further 96 build two graph convolutional modules on both the Teacher 97 Network and the Student Network to reason the relationship 98 among different people. Besides, we propose a new approach 99 for segmenting group activities in untrimmed videos, which is 100 based on the intermediate features of our model and temporal 101 convolutional networks [12]. We evaluate our approach on 102 the Volleyball dataset, Collective Activity Dataset, Collective 103 Activity Extended Dataset and Choi's Dataset, where the 104 experimental results show that the SPTS networks outperform 105 the state-of-the-arts for group activity analysis. 106

107

138

149

153

Our main contributions are summarized as follows:

- 1) In contrast to recent works for group activity recognition 108 which utilize the appearance clues only, we have devel-109 oped a Teacher Network to leverage the prior knowl-110 edge in the semantic domain, which requires no extra 111 labelled data and a little additional computational time 112 cost. 113
- 2) Different from existing self-attention based works, we 114 have explored the discriminative information of different 115 people by transferring the semantics-preserving attention 116 learned by the Teacher Network to the Student Network 117 in the appearance domain. Towards this, we equip the 118 Teacher Network and Student Network with two attention 119 modules and design an objective function which enforces 120 the Student Network to mimic the Teacher Network. 121 To our best knowledge, these are original efforts lever-122 aging attention in both semantics and appearance clues, 123 to perform group activity recognition. 124
- Unlike most conventional works which model the features 3) 125 of people as regular tensors, we have constructed two 126 types of graph for different people according to their 127 spatial coordinates, and built two graph convolutional 128 modules on the Teacher Network and Student Network to 129 reason about the relationship of different people. Exten-130 sive experimental results on four widely used datasets 131 have shown the effectiveness of our proposed method. 132
- We have extended our method for action segmentation 133 task based on its intermediate features. With the new 134 designed model, the temporal intervals of group activities 135 in an untrimmed sequence can be accurately segmented 136 and our method achieves very competitive performance 137 on this task.

It is to be noted that a preliminary conference version of 139 this work was initially presented in [13]. As an extension, our 140 SPTS with two new graph convolutional modules can better 141 exploit the interaction information of different people. More-142 over, we have conducted experiments on other two datasets 143 and provided more in-depth analysis on the experimental 144 results. Furthermore, we have extended our approach on action 145 segmentation task for untrimmed videos and demonstrate its 146 effectiveness. Besides, we have presented analysis on the 147 computational time cost of our work. 148

# II. RELATED WORK

In this section, we briefly review four related topics: 150 1) group activity recognition, 2) attention-based models, 151 3) knowledge distillation, and 4) graph convolutional network. 152

# A. Group Activity Recognition

Activity recognition is one of the most important issues in 154 computer vision [14]–[18], where group activity recognition is 155 an active sub-topic and various methods have been explored 156 in recent years [1]–[7], [19]. These methods can be roughly 157 divided into two categories: hand-crafted feature based and 158 deep learning feature based methods. For the first category, 159 a number of researchers fed hand-crafted features into graph-160 ical models to capture the structure of group activity. For 161 example, Lan et al. [9] presented a latent variable framework 162

to model the contextual information of person-person inter-163 action and group-person interaction. Hajimirsadeghi et al. [1] 164 developed a multi-instance model to count the instances in a 165 video for group activity recognition. Shu et al. [10] employed 166 AND-OR graph formalism to jointly group people, recognize 167 event and infer human roles in aerial videos. However, these 168 methods relied on hand-crafted features, which require strong 169 prior knowledge and were short of discriminative power to 170 capture the temporal cue. 171

For the deep learning based methods, numbers of works 172 have been proposed to leverage the discriminative power 173 of deep neural network for group activity recognition. For 174 example, Ibrahim et al. [2] proposed a hierarchical model 175 with two LSTM networks, where the first LSTM captured 176 the dynamic cues of each individual person, and the second 177 LSTM learned the information of group activity. Shu et al. [3] 178 extended this work by replacing the softmax layer of the RNN 179 with a new energy layer to improve reliability and numerical 180 stability of inference. Wang et al. [6] built another LSTM 181 network upon this work to capture the interaction context of 182 different people. More recently, Ibrahim et al. [20] developed 183 a Hierarchical Relational Network architecture to calculate the 184 relational representation of people and describe their potential 185 interactions. However, the works mentioned above mainly 186 focused on the appearance domain, which ignored the semantic 187 relationship between the individual actions and group activity. 188 More recently, Li et al. [4] presented a SBGAR scheme, which 189 generated the captions of each video and predicted the final 190 activity label based on these captions. However, the generated 191 captions were not always reliable, and the inferior captions 192 will do harm to the final process of recognition. To this 193 end, we simultaneously explore the contextual relationship of 194 individual actions and group activity in both semantic and 195 appearance domains, and employ the semantic knowledge to 196 enhance the performance of vision task. 197

# 198 B. Attention-Based Models

Attention-based model is motivated by the attention mech-199 anism of primate visual system [21], [22]. It aims to select 200 the most informative parts from the global field. In the past 201 two decades, attention-based models have been widely applied 202 203 into the realm of natural language processing (e.g., machine translation [23], [24]), computer vision (e.g., video face recog-204 nition [25], [26], person re-identification [27], object local-205 ization [28]), and their intersection (e.g., image caption [29], 206 video caption [30] and visual question answering [31]). 207 As for human action/activity recognition, Liu et al. [32] 208 developed global context-aware attention LSTM networks to 209 select the informative joints in skeleton-based videos. Further-210 more, Song et al. [33] proposed a spatial-temporal attention-211 based model to learn the importance of different joints and 212 different frames. Different from these two works [32], [33], 213 we employ the attention model to allocate different weights to 214 different people in a group for RGB-based activity recognition. 215 Although a few works [5], [11] have exploited attention-216 based models for group activity recognition, they only 217 applied "self-attention" scheme and were incapable to explain 218 the physical meaning of the learned attention explicitly. 219

Different from these methods, our SPTS networks distill the attention knowledge in the semantic domain to guide the appearance domain, which utilize the semantic information adequately and make the learned attention interpretable by further showing the visualization results.

# C. Knowledge Distillation

The concept of "knowledge distillation" is originated from 226 the work [34] by Hinton *et al.*, which aims to transfer the 227 knowledge in a "teacher" network with larger architecture 228 and higher performance to a smaller "student" network. They 229 enforced a constraint on the softmax outputs of the two net-230 works when optimizing the student network. After that, several 231 works have been proposed to regularize the two networks 232 based on the intermediate layers [29], [35], [36]. For example, 233 Yim et al. [36] utilized flow of solution procedure (FSP) 234 matrix, which were generated based on feature maps of two 235 layers, to transfer knowledge in teacher network to student 236 network. Chen et al. [37] employed technique of function-237 preserving transformations to accelerate the learning process 238 of student network. The most related work to ours is [29], 239 which also utilized the information across the attention mod-240 ules of two networks. Different from [29], where the inputs 241 of the two networks were both images and the networks 242 architecture were similar, our work explores the knowledge 243 in two different domains (semantic domain and appearance 244 domain) and utilizes the additional recurrent neural network to 245 address a more challenging task of group activity recognition. 246

# D. Graph Convolutional Network

Recently, there has been progress in the formulation of 248 convolutional neural network on graphs (i.e. graph convolu-249 tional network) [38]–[41] thanks to the development of graph 250 signal processing (GSP) [42]. Given inputs on the nodes of the 251 graph, the graph convolutional network (GCN) aims to learn 252 representative features like standard CNN, which sheds lights 253 on new possibilities to adopt data-driven method and perform 254 convolutional operator on non-Euclidean space. Computer 255 vision has also benefited from GCN in recent years [43], [44]. 256 For example, Wang et al. [45] considered the semantic 257 embeddings as different nodes of the knowledge graph, and 258 adopted graph convolutional network to promote the problem 259 of zero-shot recognition. Wang et al. [46] proposed a Graph 260 Reasoning Model (GRM) to study the problem of social 261 relationship understanding. For human action recognition, 262 several works [47]-[49] have been proposed to develop graph 263 convolutional network for skeleton-based action recognition. 264 Unlike these works which regarded the coordinates of human 265 joints as the nodes of the graph, we construct the nodes of the 266 graph according to the features of individual person in both 267 semantic domain and appearance domain. Then, we employ 268 two graph convolutional modules to model the relationship of 269 different people and enhance the recognition performance. 270

# III. Approach

The motivation of this work is to adequately explore the <sup>272</sup> information in both appearance domain and semantic domain <sup>273</sup>

225

247



Fig. 2. A framework of our proposed SPTS networks, which contain two sub-networks. We first train the Teacher Network, which models relationship between words of individual actions and the word of group activity. Next, we train the Student Network, which takes a set of tracklets as input and predicts the label of group activity. We enforce three types of constraints during the training process of Student Network, *i.e.*, semantics-preserving attention constraint, knowledge distillation constraint and classification constraint.

for group activity recognition. In this section, we first formulate the problem, then we present the details of our SPTS
networks and introduce how to build several graph convolutional modules on the SPTS. Finally we discuss the difference
of our models with other related works.

# 279 A. Problem Formulation

We denote a tri-tuple (V, y, z) as a training sample for a 280 video clip, where V is the specific video and z is the ground-281 truth label for group activity. Let  $Y = \{y_n\}_{n=1}^N$  denote the labels of individual actions, where  $y_n$  represents the label 282 283 corresponding to the *n*th person. The goal of group activity 284 recognition is to infer the final label z corresponding to V285 during testing phase. Previously, researchers usually utilize 286 a set of tracklets of the people in the video as inputs. The 287 tracklets are denoted as  $X = \{x_1^t, x_2^t, ..., x_n^t, ..., x_N^t\}_{t=1}^T$  where 288 t represents the time stamp of the tth frame. We follow this 289 problem setting in our work. 290

# 291 B. SPTS Networks

Our SPTS networks consist of two subnetworks, namely 292 Student Network and Teacher Network. Fig. 2 illustrates the 293 pipeline of SPTS networks. In this framework, the Student 294 Network aims to predict the final label z given a set of 295 tracklets from an input video in the appearance domain, while 296 the Teacher Network aims to model the relationship between 297 the words of individual actions  $Y = \{y_n\}_{n=1}^N$  and the word 298 of group activity z in the semantic domain. It is reasonable 299 that Teacher Network tends to achieve comparable or better 300 performance than Student Network, because individual action 301 labels are powerful low-dimensional representations for the 302 task of group action recognition, which is also demonstrated 303 in the Experiments section. Additionally, we find the Teacher 304 Network and Student Network are complementary in classi-305 fication results, which indicates that jointly considering the 306 semantic domain and appearance domain will help. However, 307 the ground-truth individual labels  $Y = \{y_n\}_{n=1}^{\hat{N}}$  are not 308

available during the testing stage. A natural way to address this issue is to employ the knowledge of the Teacher Network to guide the training process of the Student Network. We now detail the proposed SPTS networks as follows.

1) Student Network: The goal of our Student Network is to learn a model  $z = \mathbf{S}(X; \theta_s)$  to predict the label of group activity given a set of tracklets in a video clip, where  $\theta_s$  is the set of learnable parameters of the Student Network. For a fair comparison, we utilize the off-the-shelf tracklets provided by [2], [7].

In order to capture the appearance information and temporal 319 evolution of each single person, we employ a DCNN network 320 and an LSTM network to extract features of X, which is a 321 similar scheme according to [2]. Then, we concatenate the fea-322 tures of the last fc layers of the DCNN and the LSTM network. 323 The concatenation, denoted as  $G = \{g_1^t, g_2^t, ..., g_n^t, ..., g_N^t\}_{t=1}^T$ , 324 represents the temporal feature of each individual person. 325 Sequentially, we calculate the score  $s_n^t$  which indicates the 326 importance of the *n*th person as: 327

$$s_n^t = tanh(W_1g_n^t + b_1),$$
 (1) 320

where  $W_1$  and  $b_1$  are the weighted matrix and biased term. <sup>329</sup> The activation weight we allocate to each person is obtained <sup>330</sup> as follow: <sup>331</sup>

$$\beta_n^t = exp(s_n^t) / \sum_{j=1}^N exp(s_j^t), \qquad (2) \quad {}_{332}$$

where  $\beta_n^t$  is the score normalized by a softmax function. Instead of conventional aggregation methods like max-pooling or mean-pooling, we fuse the feature of each individual person at time-step *t* as:

$$w_{agg}^{t} = \sum_{n=1}^{N} \beta_{n}^{t} \cdot g_{n}^{t} \,. \tag{3}$$

In this way, the set of activation factors  $\{\beta_n^t\}_{n=1}^N$  control the some contribution of each person to the aggregated feature  $w_{agg}^t$ .

Having obtained  $w_{agg}^t$ , the aggregated features of each frame, 340 we feed them into another group-level bidirectional LSTM 341 network. The output features are sent into an fc layer activated 342 by a softmax function to obtain the final label of the group 343 activity. 344

2) Teacher Network: As illustrated above, our Student 345 Network can be regarded as an extension of the hierarchical 346 deep temporal model [2] by adopting a typical self-attention 347 mechanism. However, in such a scheme, the labels of indi-348 vidual actions and group activities are utilized to supervise 349 the discriminative feature learning, while their corresponding 350 relationship, which captures the dependency of the individual 351 actions and group activities in the semantic domain, is rarely 352 used. In this section, we introduce a Teacher Network, which 353 aims to learn a model  $z = \mathbf{T}(Y; \theta_t)$  to integrate the labels of individual actions  $Y = \{y_n\}_{n=1}^N$  into a label of group 354 355 activity z. Note that our Teacher Network essentially addresses 356 an NLP-related task, where attention mechanism also shows 357 its advantage. Based on this, we develop our Teacher Network 358 by introducing an attention scheme, which is similar to our 359 Student Network. 360

Given a set of individual action labels  $Y = \{y_n\}_{n=1}^N$  as the 361 input of our Teacher Network, we first encode them into a sequence of one-hot vectors  $F_{oh} = \{f_{oh,n}\}_{n=1}^N$ , where  $f_{oh,n} \in \mathcal{F}_{oh}$ 362 363  $R^{C}$  and C is the number of individual action category. Then 364 we embed the  $F_{oh} \in \mathbb{R}^{P \times C}$  into a latent space as: 365

$$f_{em,n} = ReLU(W_2 f_n + b_2),$$
 (4)

where  $W_2$  and  $b_2$  are the weighted matrix and biased term, 367 *ReLU* denotes the nonlinear activation function [50]. Then 368 another attention mechanism, which is corresponding to that 369 of the Student Network, is derived as follow: 370

$$s_n = tanh(W_3 f_{em,n} + b_3),$$
 (5)

372

36

$$\alpha_n = exp(s_n) / \sum_{j=1}^N exp(s_j), \qquad (6)$$

388

 $v_{agg} = \sum_{1}^{N} \alpha_n \cdot f_{em.n} \, .$ 

Having obtained the  $v_{agg}$ , we feed it into an fc layer 374 followed by a softmax activation to predict the final label. 375 We train the Teacher Network using the ground-truth labels 376 of Y and z. It is relatively easy to classify a set of words in 377 the semantic domain, thus the Teacher Network will achieve 378 higher performance as illustrated in the Experiments section. 379 3) Semantics-Preserving Attention Learning: As we 380 described, there are two attention modules in our method 381 and they both work separately via a self-attention scheme. 382 Noticing the fact that they both model the importance of 383 different people, a valid question is why not jointly consider 384 these two modules. More specially, as the Teacher Network 385 directly takes the ground-truth label of individual actions as 386 inputs, it is reasonable that its performance is better than 387 the Student Network, which takes the tracklets as inputs and

requires a more complex feature learning process before the 389 attention module. 390 Based on this reason, we aim to use the attention knowl-391 edge of the Teacher Network to guide the Student Network. 392

# Algorithm 1 SPTS

<b>Input</b> : Training samples: $\{X, Y, z\}$ , Parameters: $\Gamma$
(iterative number) and $\epsilon$ (convergence error).
<b>Output</b> : The weights of the Student Network $\theta_s$ .
// Teacher Network Training:
Optimize the parameter $\theta_t$ of the Teacher Network w

vith (Y, z).

// Student Network Training:

Finetune the DCNN and the train first LSTM with (X, Y) [2].

Extract features G from X.

Initialize  $\theta_s$ .

Perform forward propagation.

Calculate the initial  $J_0$  by (8).

for  $i \leftarrow 1, 2, ..., \Gamma$  do Update  $\theta_s$  by back propagation through time (BPTT). Perform forward propagation.

Compute the objective function  $J_i$  using (8).

If  $|J_i - J_{i-1}| < \epsilon$ , go to **Return**.

end

(7)

**Return:** The parameters  $\theta_s$  of the Student Network.

In practice, we first train the Teacher Network  $\mathbf{T}(Y; \theta_t)$  with 393 the provided labels of training samples. Then, we enforce the 394 Student Network to absorb the teacher's knowledge during the 395 learning process via a total loss function defined as below: 396

$$J = J_{CLS} + \lambda_1 \ J_{SPA} + \lambda_2 \ J_{KD}$$
<sup>397</sup>

$$= -\sum_{l=1}^{L} \mathbb{1}(z=l)log(P_S^l)$$
390

$$+\lambda_1 \frac{1}{N} \sum_{n=1}^{N} (\alpha_n - \frac{1}{T} \sum_{t=1}^{I} \beta_n^t)^2$$
 399

$$+ \lambda_2 \|P_T - P_S\|_2^2 \tag{8}$$

Here  $\lambda_1$  and  $\lambda_2$  are the hyper-parameters to balance the 401 effects of two different terms to make a good trade-off. The 402 physically interpretations of the  $J_{CLS}$ ,  $J_{SPA}$  and  $J_{KD}$  are 403 respectively explained as below. 404

The first term  $J_{CLS}$  represents classification loss for activity 405 recognition. We calculate the categorical cross-entropy loss, 406 where  $\mathbb{1}$  is the indicator function which equals 1 when the 407 prediction z = l is true and 0 otherwise. Here l and L denote 408 the predicted label and the number of the total activity cat-409 egories. The softmax output  $P_S^l$  represents the corresponding 410 class probability of the Student Network. The second term 411  $J_{SPA}$  aims to enforce the student's attention to preserve the 412 teacher's semantics attention. We adopt the mean squared 413 distance for these two types of attention. The third term  $J_{KD}$ 414 denotes the loss of knowledge distillation [34], in which  $P_T$ 415 and  $P_S$  are the softmax outputs of the Teacher and Student 416 Network respectively. 417

To optimize (8), we employ the back propagation through 418 time (BPTT) algorithm [51] for learning all the parameters  $\theta_s$ 419 of our Student Network. We summarize the pipeline of our 420 SPTS method in Algorithm 1. Note that the Teacher Network 421

443

456

only guides the Student Network during the training phase, 422 as the ground-truth label  $Y = \{y_n\}_{n=1}^N$  is not available during 423 the testing stage. 424

#### C. SPTS + GCN425

Since a group of people can be considered as a graph-based 426 structure, where the node and edge represents each individual 427 person and the relationship between two people respectively, 428 we further build two graph-based modules upon our SPTP 429 networks to adequately explore the contextual information of 430 different people for group activity recognition. 431

1) Graph Construction: We construct a graph  $\mathcal{G}(U, A)$  to 432 model each frame, where U and A are the nodes sets and 433 adjacency matrix respectively. On the one hand, we denote 434  $U = \{u_1, u_2, ..., u_N\}$ , where  $u_n \in D$  is corresponding to the 435 feature of the *n*th person. On the other hand, motivated by 436 the fact that, the relationship of different people are highly 437 correlated to the distance among them, we define the adjacency 438 matrix A according to the spatial coordinates of different 439 people as follow: 440

$$a_{mn} = exp(-\frac{||c_m - c_n||_2^2}{2}), \qquad (9)$$

where  $c_m$  represent the central location of the *m*th person: 442

$$c_m = (\gamma \, \frac{x_{m,mid}}{W \, I}, \gamma \, \frac{y_{m,mid}}{H \, I}). \tag{10}$$

Here, W\_I and H\_I are the width and height of each frame 444 respectively.  $x_{m,mid}$  and  $y_{m,mid}$  are the central positions of the 445 input tracklets at the x axis and y axis. The  $\gamma$  is a scale factor, 446 where we set it to be 10 empirically. In this way, we embed the 447 spatial information into the adjacency matrix A. If two people 448 m and n approach each other in the space, the corresponding 449  $a_{mn}$  will have a large value, and vice versa. 450

2) Graph Convolutional Layer: Since the graph of peo-451 ple lie in a non-Euclidean space, we leverage the graph-452 based convolutional Networks (GCN) [39] to learn the spatial 453 dependency between different people. Mathematically, we can 454 represent a layer of the graph convolution as: 455

$$Z = AUW, \tag{11}$$

where W are the learned parameters. Unlike conventional 457 convolutional operator that reasons about the regular structure 458 locally, the graph convolutional layer passes messages among 459 different nodes and updates each nodes according to the pre-460 defined adjacency matrix A, which allows us to better capture 461 the contextual information among different people. Moreover, 462 we can stack multiple layers of graph convolution to better 463 model the non-linear structure among people. 464

3) Building GCN Upon SPTS: Fig. 3 displays the illus-465 tration of building GCN upon our SPTS. For the Teacher 466 Network, we perform graph convolution on the one-hot vector 467  $F_{oh}$  of each video clip: 468

$$Z_{teacher} = AF_{oh}W_{teacher}, \qquad (12)$$

where A is obtained based on the middle frame of the video 470 clip. The output feature  $Z_{teacher}$  is then fed into the attention 471 mechanism of the Teacher Network. 472



Fig. 3. Flowchart of building graph convolutional modules upon the SPTS networks. We develop two graph convolutional modules for better exploring the contextual information of different people. We construct two types of graph according to the spatial coordinates of different people. The graph for the Teacher Network is built based on the one-hot encoding vector  $F_{oh}$ , while the graph for the Student Network is constructed according to the extracted feature G from the input tracklets. The two graphs are sent into two graph convolutional modules to pass messages of different nodes. The output features are then fed into the two attention modules of the SPTS networks, respectively

For the Student Network, we feed  $G^t = \{g_1^t, g_2^t, ..., g_N^t\},\$ 473 the features of N people at the time stamp t, into the graph convolutional layer:

$$Z_{student}^{t} = A^{t} G^{t} W_{student}, \qquad (13) \quad {}^{476}$$

where  $A^t$  is calculated based on the tracklets of the *t*th frame. 477 We also perform instance-normalization [52] and non-linear 478 activation (ReLU) on the output feature  $Z_{student}^{t}$  before it is sent into the next layer. We stack three graph convolutional 479 480 layers for the Student Network, as the input  $G^t$  lies in a high-481 dimension space. The  $G^t$  at different time stamps t share the 482 same parameter  $W_{student}$ , we concatenate  $Z_{student}^{t}$  from 1 to T as  $Z_{student} = (Z_{student}^{1}, ..., Z_{student}^{T})$ , and then sent  $Z_{student}$ 483 484 into the attention module of the Student Network. The effects 485 of the number of graph convolutional layer will be explored 486 in the Experiments section. 487

# D. Discussions

We discuss the difference of our methods with other two categories of DNN-based methods in this subsection.

The first category, such as HTDM [2] and its variants [3] 491 shown in Fig. 4(a), mainly focus on the appearance domain. 492 They first learn features of individual person with an LSTM 493 network, then aggregate them into group representations with 494 a function  $f_1$ , and finally recognize the activity based on the 495 group representations with another LSTM network. The labels 496 of individual actions Y and group activity z were respectively 497 used to supervise the training process of the first and second 498 LSTM networks. But the corresponding relationship of Y and z 499 have not been utilized explicitly. Moreover, the function 500  $f_1$  turned to be max-pooling or mean-pooling, which lacks 501 physical meaning. 502

The second category, such as SBGAR [4] displayed 503 in Fig. 4(b), focuses on the semantic domain. This method 504

474 475

488

489



Fig. 4. Comparison of different DNN-based frameworks for group activity recognition. The solid lines, dashed lines and green arrow denote the process of forward propagation, backward propagation and semantics-preserving attention learning respectively. Method in (a) first extracts features of individual action, then aggregates them into group representations with  $f_1$ , and finally recognizes the activity based on the group representations. Approach in (b) first generates captions (*i.e.*, individual action labels) of video frames, and recognizes the activity based on these captions by  $f_2$ . Our method in (c) first employs two graph convolutional modules to capture the contextual information of features in both semantic and appearance domain. Then we learn  $f_4$  to classify the group activity label based on the learned features in the semantic domain. Finally, we employ the attention knowledge in  $f_4$  to guide  $f_3$  when aggregating features in the appearance domain to make the final prediction.

directly generates the caption to describe the video frames, 505 and utilizes the captions to classify the group activity with a 506 function  $f_2$ . The individual actions Y were used to supervise 507 the process of caption generation and the group activity z508 was utilized to supervise the learning process of  $f_2$ . However, 509 as the group label is sensitive to the captions, the inaccurate 510 generated captions will do harm to the final recognition results. 511 Different from these methods, our approach in Fig. 4(c), 512 adequately leverage the information in the appearance domain 513

and the semantic domain for group activity recognition. 514 We distill the knowledge in  $f_4$  learned in the semantic domain 515 to guide the training process of  $f_3$  in the appearance domain. 516 Moreover, we have employed two graph convolutional mod-517 ules to further reason the dependency of different people and 518 enhanced the final recognition performance. 519

### E. Exploration on Temporal Segmentation for Group Activity 520

Temporal segmentation (a.k.a. action segmentation) aims 521 to segment actions in untrimmed videos and recognize their 522 action labels. Although it has attracted growing attention 523 in recent years [12], [53]-[56], few attempts on temporal 524 segmentation for group activity have been devoted due to the 525 scarcity of annotated datasets and complicated relationship of 526 different people. In order to see how our method performs on 527 this task, we have made explorations as follows. 528

Fig. 5 presents the illustration of incorporating our method 529 with temporal convolutional networks (TCN) [12] for group 530 activity segmentation. Since our method takes the tracklets 531 of N people in T frames as input, we first divide the input 532 video into L clips and the length of each clip is T frames. 533 Then we employ faster-RCNN [57] to detect people in each 534 frames, and align the cropped people in T frames according 535 to their locations. Through this pre-process, we obtain a set 536 of tracklets and choose N of them according to the top-N 537 detection scores in the first frames of the clip. Then we adopt 538 a DCNN and LSTM network to extract the features  $\{F_1^l\}_{l=1}^L$ 539



Fig. 5. Flowchart of combining our method with temporal convolutional networks (TCN) [12] for group activity segmentation. The input of the approach is an untrimmed video with  $L_{total}$  ( $L_{total} = 1000$ ) frames, we first divide it into L (L=100) clips and the length of each clip is T (T =10) frames. Then we generate the tracklets based on the mask-rcnn detector and the locations of different people. Similar with the trimmed setting, the tracklets are feed into a DCNN and LSTM network to extract features of individual actions. The extracted features are sent into our model (SPTS Network + GCN) and generate the features of group activities for each clips. Finally, we concatenate these clip-based features to a video-based feature and utilize TCN model to learn the segmentation results. For the 1-th clips, the  $F_1^l$  and  $F_2^l$  are corresponding to the G and  $\{w_{agg}^t\}_{t=1}^T$  in Fig.2.

of the input tracklets, where  $F_1^l$  is a tensor with the shape of 540  $N \times T \times d$ . Here d is the summed dimension of the last fc layers 541 in the DCNN and LSTM networks. The features of individual 542 actions are fed into our model (SPTS Network + GCN). 543 Finally, we concatenate the output features  $\{F_2^l\}_{l=1}^L$  into a video-based feature  $F_3 = concat(F_2^1, F_2^2, ..., F_2^L)$  and sent 544 545 it into the TCN model to obtain the segmentation results. 546

# **IV. EXPERIMENTS**

In this section, we conducted experiments on three 548 datasets for group activity recognition, including volleyball 549 dataset [59], collective activity (CA) dataset [60] and collective 550 activity extended (CAE) dataset [8]. The experimental results and analysis are described in details as follows. 552

551



Fig. 6. Examples of the pair-wise representative frames from three different datasets we used. For each group, the RGB-based pictures are presented on the left, while the corresponding optical flows extracted by Flownet 2.0 [58] are shown on the right. (a) Volleyball dataset. (b) Collective activity dataset. (c) Collective activity extended dataset. (d) Choi's dataset.

# 553 A. Datasets and Experiment Settings

1) Volleyball Dataset [59]: The Volleyball dataset is cur-554 rently the largest dataset for group activity recognition. It con-555 tains 55 vollevball videos with 4830 annotated frames. There 556 are 9 individual action labels (waiting, setting, digging, falling, 557 spiking, blocking, jumping, moving and standing) and 8 group 558 activity categories (right set, right spike, right pass, right 559 winpoint, left winpoint, left pass, left spike and left set) in 560 this dataset. We employ the evaluation protocol in [59] to 561 separate the training/testing sets. We employ the metrics of 562 Multi-class Classification Accuracy (MCA) and Mean Per 563 Class Accuracy (MPCA) on this dataset. 564

2) Collective Activity (CA) Dataset [60]: The Collective 565 Activity Dataset is a widely used benchmark for the task of 566 group activity recognition. It comprises 44 video clips, anno-567 tated with 6 individual action classes (NA, crossing, walking, 568 waiting, talking and queueing) and 5 group activity labels 569 (crossing, walking, waiting, talking and queueing). There are 570 also 8 pairwise interaction labels, which we do not utilize in 571 this paper. We split the training and testing sets following the 572 experimental setup in [9]. 573

As suggested in [60] that originally presented the dataset, the "walking" activity is rather an individual action than a collective activity. To address this, we follow the experimental setup in [6], to merge the class of "walking" and "crossing" as a new class of "moving". We report the Mean Per Class Accuracy (MPCA) of the four activities on the CA dataset, which can better evaluate the performance of the classifiers.

*3) Collective Activity Extended (CAE) Dataset [8]:* The Collective Activity Extended Dataset contains 7 individual action labels and 6 group activities categories. It replaces the "walking" activity with other two activities of "dancing" and "jogging" in the CA Dataset. We adopted the training and testing splits used in [61] to train our models.

4) Choi's Dataset [7]: The Choi's dataset comprises 587 32 videos, which are annotated with 3 individual actions 588 (walking, standing still, and running), and 6 group activi-589 ties (gathering, talking, dismissal, walking together, chasing, 590 and queueing). The dataset also provided 8 pose labels and 591 9 interaction labels which we did not utilize. We followed the 592 standard experimental protocol of the 3-fold cross validation, 593 which was adopted in [7]. 594

595 5) Untrimmed Volleyball Dataset [59]: The untrimmed 596 Volleyball dataset consists of 54 long videos of Volleyball 597 datasets,<sup>1</sup> which is for temporal segmentation. The duration

of each video varies from 76.76 minutes to 185.13 minutes. 598 Since the length of these videos are too long for analysis 599 and only numbers of temporal intervals have been annotated 600 in [2]. We proceed them in to 837 clips according to the 601 annotation [2], where each clips has 1000 frames. We chose 602 this length as it is comparable with the duration of video clips 603 in GTEA dataset [62] and 50 Salads dataset [63] evaluated 604 by TCN [12]. We finally obtained 612 clips for training and 605 225 clips for testing. There are 8 group activity labels (the 606 same with [2]) and a background label. We report the F1 score 607 at frame level, which is computed as: 608

$$F1 = \frac{2 \times precision \times recall}{precision + recall}.$$
 (14) 601

610

# B. Implementation Details and Baselines

 Group Activity Recognition: Our proposed methods were built on the Pytorch toolbox and implemented on a system with the Intel(R) Xeon(R) E5-2660 v4 CPU @ 2.00Ghz. We trained our model with two Nvidia GTX 1080 Ti GPUs and tested it with one GPU.

For the Teacher Network, we took the ground-truth label of 616 each individual action as input, and the one-hot vectors were 617 projected through an fc layer. The embedded features were 618 weighted and summed based on different weights learned by 619 the self-attention mechanism, which indicates the importance 620 of different people. The aggregated features were then fed into 621 an fc layer for classification. The Teacher Network was trained 622 with the Adam optimization method with 16 as the batch size. 623 And the initial learning rate was 0.003. 624

For the Student Network, we first finetuned VGG net-625 work [64] pretrained on ImageNet [65] to extract CNN fea-626 tures of the tracklets. The features of the last fc layer were 627 fed into a LSTM network with 3000 nodes. The concatenated 628 features of VGG and LSTM networks were then fed into an 629 fc layer with the size of 512 to cut down the dimension. The 630 importance of each person on each frame was generated by 631 the attention mechanism, and the embedded features of each 632 person were then summed by weight. The weighted features 633 were then fed into a bidirectional LSTM network with the 634 hidden size of 128. The output features were fed into an fc 635 layer for classification. During the Teacher guided training 636 process, the Student Network was optimized with Adam and 637 the initial learning rate was 0.00003. As for ratio of different 638 parts of losses, we set  $\lambda_1 = \lambda_2 = 1$ . The batch size was set 639 to be 16. 640

In order to better explore the motion information of the video and inspired by the success of two-stream network architecture [18], we computed the optical flow between two adjacent video frames using Flownet 2.0 [58]. We extracted

<sup>&</sup>lt;sup>1</sup>The original volleyball dataset provided trimmed clips and the names of 55 long videos. However, the 21-th video cannot be found according to its names. Moreover, due to the changes of frame rate on YouTube, 8 videos are incorrectly aligned with the temporal annotation provided in [2]. To address this, we spent 2 days refining the annotations to ensure their correctness.

the DCNN and LSTM features of optical flow tracklets, and
concatenated them with the features of the original RGB
tracklets before the attention module of the Student Network.
We report the performance of the following baseline methods and different versions of our approach:

649 • HDTM [2]: A hierarchical framework with two LSTM 650 models. The first LSTM network took the features 651 extracted from the tracklets of each person as input, and 652 was trained with the supervision of the individual action 653 label. The input of the second LSTM network was the 654 aggregation of features learned by the first LSTM, and 655 was trained with the supervision of the group activity 656 label. 657

- Ours-teacher\*: The Teacher Network directly took the 658 ground-truth labels of the individual actions as input 659 during both training and testing phases. Hence, it is 660 not fair to directly compare the performance of Teacher 661 Network with other methods, which are inaccessible to 662 the ground-truth labels of the individual actions during 663 testing phase. We report the performance of Ours-teacher\* 664 only for reference. 665
- Ours-teacher: During the training phase, we used the ground-truth label of each individual action as input to train the Teacher Network. During the testing stage, we used the individual action label learned from the first LSTM of HDTM to predict the final group activity label.
  - Ours<sub>-SA</sub> (self-attention): An original model of our Student Network, which can be regarded as adding a self-attention module upon the HDTM [2].

671

672

673

677

678

679

680

681

682

- Ours\_*SPA* (semantics-preserving attention): A version of model which employed the attention knowledge in Teacher Network to help the training of Student Network.
  - Ours<sub>-SPA+KD</sub> (knowledge distillation): A model of combining the knowledge distillation loss [34] with Ours<sub>-SPA</sub>.
  - Ours<sup>†</sup>-x: Models of combining the optical flow input based on the original Ours-x.
  - Ours-teacher\* + GCN: Building the graph convolutional module upon the Teacher Network.

• Ours+GCN<sub>-SA</sub>, Ours+GCN<sub>-SPA+KD</sub>, Ours<sup>†</sup> +GCN<sub>-SA</sub> and Ours<sup>†</sup> +GCN<sub>-</sub> SPA+KD: Models of equipping the graph convolutional module with Ours<sub>-SA</sub>, Ours<sub>-SPA+KD</sub>, Ours<sup>†</sup><sub>-SA</sub> and Ours<sup>†</sup><sub>-</sub> SPA+KD.

2) Temporal Segmentation for Group Activity: During 687 experiments, we first pretrained our model on the trimmed 688 Volleyball dataset, and finetuned it on the untrimmed dataset 689 to extract features. We report the segmentation results of 690 comparing methods in two categories: image-level methods 691 and person-level methods. The first category consists of two 692 methods, which took the whole images as input directly: 693 (1) VGG16 [64]: We employed VGG16 network pretrained 694 on ImageNet [65], and finetuned it on the training set 695 of untrimmed Volleyball to predict the frame-level labels. 696 (2) TCN [12]: We used the features of the fc7 layer in 697 VGG16 to train the TCN models. The second category com-698 prises three approaches, which were based on the tracklets of 699 different persons: TCN<sub>-SA</sub>, TCN<sub>-SPA+KD</sub>, TCN-GCN<sub>-SPA+KD</sub>. 700 They denote using the methods  $Ours_{-SA}$ ,  $Ours_{-SPA+KD}$ , 701 Ours-GCN\_SPA+KD for feature extraction respectively. 702

TABLE I

COMPARISON OF THE GROUP ACTIVITY RECOGNITION ACCURACY (%) ON THE VOLLEYBALL DATASET.<sup>†</sup> DENOTES THAT THE MODEL TAKES BOTH RGB IMAGES AND OPTICAL FLOWS AS INPUTS

Method	MCA	MPCA
CERN-2 [3]	83.3	83.6
SSU [5]	89.9	_
SRNN [66]	83.5	_
RCRG [20]	89.5	_
Ours-teacher*	88.3	84.4
Ours-teacher* + GCN	92.3	90.7
Ours-teacher	69.3	66.8
Baseline-HDTM [2]	86.8	85.8
Ours – SA	87.1	86.1
Ours _ SPA	89.3	89.2
Ours $_{-SPA + KD}$	89.3	89.0
$Ours^{\dagger} - SA$	87.7	87.0
Ours <sup>†</sup> – SPA	89.6	89.5
$Ours^{\dagger} - SPA + KD$	90.7	90.0
Ours + GCN $_{-SA}$	89.2	88.8
Ours + GCN $_{-SPA + KD}$	90.4	89.3
$Ours^{\dagger} + GCN_{-SA}$	90.4	90.5
$Ours^{\dagger} + GCN - SPA + KD$	91.2	91.4

# C. Results on the Volleyball Dataset

We first evaluate our proposed methods on the Volleyball 704 dataset. We follow [2] to separate players into two groups 705 on the left and right, and extend the individual action labels 706 to 18 categories (*e.g.*, "left standing", "right waiting", etc.) 707 according to their spatial coordinates. 708

1) Comparison With the State-of-the-Arts: Table I presents709the comparison performance with different approaches.710We observe that our final model (Ours<sup>†</sup> + GCN\_SPA + KD)711achieves 91.2% MCA and 91.4% MPCA, outperforming existing state-of-the-art methods for group activity recognition.713

2) Analysis on the SPTS Networks: Here we analyze 714 our semantics-preserving learning scheme. Compared with 715 the 0.3% (MCA and MPCA) improvement by the self-716 attention scheme over the baseline method, our attention-717 guided approach achieves 2.5% (MCA) and 3.2% (MPCA) 718 improvement, which demonstrates the effectiveness of our 719 proposed method. We also discover that, combining with the 720 optical flow can lead to a slight improvement on this dataset. 721 While besides, Our-teacher\*, which takes the ground-truth 722 of individual actions as the testing inputs of the Teacher 723 Network, reaches performance of 88.3% MCA, Our-teacher, 724 which utilizes the predicted individual actions as the testing 725 inputs, only attains 69.3% MCA. This is because, the Teacher 726 Network is sensitive to the inputs and the incorrected predicted 727 individual actions will greatly harm the performance of the 728 final recognition. 729

We also show several visualization results of the learned 730 attention in Fig. 7. The group activity label of Fig. 7(a) 731 is "left spike". For the self-attention model of the Student 732 Network, the model most likely focuses on those people 733 wearing different clothes in a group, e.g., the white per-734 son (SA:60) in the black team, and the yellow person (SA:62) 735 in the white team. However, these people are not exactly key 736 people for recognizing the group activity. When we employ 737 the attention model of Teacher Network, we can focus on 738 those words, which are essentially important in the semantic 739



Fig. 7. Visualization of the learned attention on the Volleyball dataset. In (a)(b)(c), for each video clip, we show the representative frame on the left, while the cropped people are shown on the right. In each dash box, we display the labels of individual actions and three types of attention score: T (Teacher Network), SA (Student Network with self-attention scheme) and SPA (Student Network with semantics-preserving attention method). The SA and SPA scores in (a)(b)(c) are averaged scores over a clips (10 frames). In (d)(e), we present the attention scores and the corresponding people in temporal domain.

domain, e.g., the spiking (T:80), and the blocking (T:51). 740 And after employing our SPTS networks, we will transfer 741 this attention knowledge from the semantic domain to the 742 appearance domain, and guide the Student Network to focus 743 on the "left spiking" person (SPA:62), who contributes most 744 745 to recognizing the final activity. The group activity label of Fig. 7(b) is "left winpoint", where there is no special people for 746 recognizing this activity. However, the self-attention scheme 747 assign the highest score to the yellow person (SA:72), which 748 does not carry key information. After employing the SPTS 749 networks, the score of this person is decreased to 47, and 750 extra attention is allocated to other people. Fig. 7(c) illustrates 751 similar results to Fig. 7(a). 752

We further present the learned attention scores on temporal domain in Fig. 7(d) and Fig. 7(e). For the "spiking" people in volleyball dataset, our SPA scores (blue ones) go up to climaxes when the players wave their hands to spike the ball, which assigns more attention to the discriminative frames.

3) Analysis on the Graph Convolutional Modules: 758 As shown in Table I, when applying the graph convolu-759 tional modules, the Teacher Network achieves 4.0% and 760 6.3% improvement on the MCA and MPCA metrics respec-761 tively. For the Student Network, Ours  $Ours^{\dagger} + GCN_{-SA}$  and 762  $Ours^{\dagger} + GCN_{-SPA + KD}$  attain 2.7% and 0.5% improvement 763 on MCA, and 3.5% and 1.4% improvements on MPCA, 764 which consistently demonstrates the effectiveness of the graph 765 convolutional modules. 766

Moreover, we have conducted experiments on adopting 767 different layers for the Teacher Network and Student Network. 768 As presented in Table II, the peaks of the Teacher Network and 769

TABLE II Comparison of the Group Activity Recognition Accuracy (%) of Different Number of Graph Convolutional Layers on the Volleyball Dataset

Number of Graph Convolutional Layers	1	3	5	7
Ours-teacher* + GCN (semantic domain)	92.3	91.3	90.9	90.4
$Ours^{\dagger} + GCN_{-SA}$ (appearance domain)	89.6	90.4	90.3	90.2

## TABLE III

Comparison of the Group Activity Recognition Accuracy (%) on the CA Dataset.  $^\dagger$  Is Defined in the Caption of Table I

Method	MPCA
Cardinality kernel [1]	88.3
CERN-2 [3]	88.3
RMIC [6]	89.4
SBGAR [4]	89.9
MTCAR [7]	90.8
Ours-teacher*	97.6
Ours-teacher* + GCN	97.6
Ours-teacher	88.2
baseline-HDTM [2]	89.7
Ours – SA	91.5
Ours – SPA	92.3
Ours – SPA + KD	92.5
Ours <sup>†</sup> – SA	94.3
Ours <sup>†</sup> – <i>SPA</i>	95.6
$Ours^{\dagger} - SPA + KD$	95.7
Ours + GCN $_{-SA}$	91.8
Ours + GCN $_{-SPA + KD}$	92.9
$Ours^{\dagger} + GCN_{-SA}$	95.4
$Ours^{\dagger} + GCN_{-SPA + KD}$	95.8

Student Network appear at one layer and three layers respectively. This is because, the dimension of input feature to the
Teacher Network is relatively low and one graph convolutional
layer is proper. For the Student Network, the dimension of
input feature is much higher, thus deeper structure is needed
to achieve a better result.

# 776 D. Results on the CA Dataset

1) Comparison With the State-of-the-Arts: Table III shows 777 the comparison with different methods on the CA dataset. 778 The MPCA results of other approaches are computed based 779 on the original confusion matrices in [1]-[4], [6], [7]. 780 We observe that, our final model ( $Ours^{\dagger} + GCN_{-SPA + KD}$ ) 781 achieves 95.8% MPCA, outperforming the state-of-the-art [7] 782 by 5.0%. Moreover, our method have improved the baseline 783 method HDTM [2] by 6.0%. Fig. 9 presents the confusion 784 matrices of the baseline methods and our SPTS networks. It is 785 clear that SPTS networks attain superior results, especially 786 for distinguishing the activity of "moving" and "waiting". 787 Besides, compared with SBGAR and Ours-teacher, which 788 directly utilized the semantic information to predict the final 789 labels, our method achieves 5.9% and 7.6% improvement, 790 which demonstrates its effectiveness. Objectively speaking, 791 we should own the major contribution to the combination 792 of the optical flow, which explicitly captures the motion 793 information of the scene. Based on this, our two semantics-794 preserving learning method and graph convolutional module 795 have further enhanced the recognition performance, which will 796 be discussed as follow. 797

Analysis on the SPTS Networks: From Table III,
 our attention-guided method brings 1.0%, 1.4% and 0.4%

improvements on the self-attention scheme of  $Ours_{-SA}$ , 800  $Ours_{-SA}^{\dagger}$  and  $Ours+GCN_{-SA}^{\dagger}$ . We notice that these improvements are less significant than those on the Volleyball dataset. This is because the setting of the CA dataset is to assign what the major people are doing to the label of group activity. Hence, attention model is not so important.

We also show the visualization of the learned attention 806 in Fig. 8. As shown in Fig. 8(a), the group activity label is 807 "waiting", hence the Teacher Network allocates more attention 808 to the words "waiting" (29) and less attention to the word 809 "moving". Guided by this information, the Student Network 810 decreases the attention (from 22 to 17) of the "moving" 811 person, which can be regarded as a noise for recognizing 812 the group activity. For Fig. 8(b), the group activity is "mov-813 ing", and it is reasonable that the Teacher Network allo-814 cates averaged score to the three individual words "moving". 815 Taught by this attention knowledge, the Student Network 816 increases the attention of the top person from 20 to 27, and 817 decreases the attention of the right person from 43 to 37, 818 so that the information of three people can be utilized 819 equally. 820

The temporal attention scores are shown in Fig. 8(c) and 821 Fig. 8(d). For the "spiking" people in volleyball dataset, 822 our SPA scores (blue ones) go up to climaxes when the 823 players wave their hands to spike the ball, which assigns 824 more attention to the discriminative frames. For the "waiting" 825 and "moving" people in CA dataset, the learned SPA scores 826 vary little over time because there is no part of particular 827 significance during these actions. 828

3) Analysis on the Graph Convolutional Modules: When 829 we apply graph convolutional modules to the SPTS networks, 830 the MPCA increases 1.1% and 0.1% over  $Ours_{SA}^{\dagger}$  and 831  $Ours_{-SPA + KD}^{\dagger}$  respectively, which also shows its effectiveness. 832 However, we observe that the improvements are not novel as 833 the results on the volleyball dataset. The reason is that the 834 volleyball dataset is the currently largest dataset for group 835 activity recognition, while the CA dataset is relatively small. 836 Since the graph convolutional module is a data-driven model, 837 more training data can bring more benefits. 838

# E. Results on the CAE Dataset

We further conducted experiments on the CAE dataset. 840 Table IV presents the comparison with different methods, 841 where our final model reaches a performance of 98.1%, 842 outperforming the existing state-of-the-art methods. The self-843 attention scheme achieves 95.0% and 95.9% recognition 844 accuracy on the RGB inputs and combining optical flows 845 respectively, where we obtains 0.9% and 1.7% improve-846 ments when applying our SPTS network. Moreover, Ours-847 teacher\* +GCN, Ours<sup>†</sup> +GCN<sub>-SA</sub> and Ours<sup>†</sup> +GCN<sub>-SPA</sub> + KD848 obtained 1.3%, 0.9% and 0.5% improvements benefiting from 849 the graph convolutional modules, which further shows the 850 effectiveness of the proposed approaches. 851

Fig. 9 presents the comparison of confusion matrices on the baseline method and our final model. For the baseline method, "waiting" is sometimes confused with the activity "crossing", and "dancing" is likely to be misclassified as "jogging". When applying our method, we clearly show the



Fig. 8. Visualization of the learned attention on the CA dataset. The definitions of T, SA and SPA are the same with those in Fig. 8.

TABLE IV Comparison of the Group Activity Recognition Accuracy (%) on the Collective Activity Extended Dataset. † Is Defined in the Caption of Table I

Method	Accuracy
CRF+CNN [61]	86.8
Structural SVM + CNN [61]	87.3
Structure Inference Machines [61]	90.2
Image Classification Model [11]	92.3
Person Classification Model [11]	95.1
Latent Embeddings Model [11]	97.9
Ours-teacher*	97.8
Ours-teacher* + GCN	99.1
Ours-teacher	96.0
baseline-HDTM [2]	94.2
Ours – SA	95.0
Ours – SPA	95.8
Ours $-SPA + KD$	95.9
Ours <sup>†</sup> – SA	95.9
$Ours^{\dagger} - SPA$	97.2
$Ours^{\dagger} - SPA + KD$	97.6
Ours + GCN $-SA$	95.6
Ours + GCN $-SPA + KD$	96.2
$Ours^{\dagger} + GCN_{-SA}$	96.8
$Ours^{\dagger} + GCN - SPA + KD$	98.1

advantages on discriminating these activities and obtain the promising recognition results.

# 859 F. Results on the Choi's Dataset

Table V presents the experimental results. In this dataset, our final model Ours<sup>†</sup> + GCN<sub>-SPA</sub> + *KD* achieves 78.1% accuracy, which is comparable with existing methods [2], [7], [60]. Objectively speaking, the performance of our method is not novel as those in the volleyball [59], CA [60] and CAE [8] datasets, and the reasons are two folds: (1) The 865 methods [7], [60] utilize the pose labels and interaction labels, 866 which are not used in our methods. (2) Our methods are data-867 driven based, while the methods [7], [60] use hand-crafted 868 features. So they have more advantages on the Choi's dataset, 869 which is the smallest compared with the other three datasets. 870 Besides, we observe that combining optical flow can bring 871 a large improvement in this dataset. This is because the 872 individual action labels of this dataset are "walking", "standing 873 still", and "running", so the features obtained with the input of 874 optical flow have much more discriminative power. Moreover, 875 we find the GCN and semantics-preserving attention scheme 876 can further lead to improvements, which demonstrates the 877 effectiveness of our proposed approaches. 878

# G. Results on the Untrimmed Volleyball Dataset

We evaluate our method for action segmentation on this 880 dataset and Table VI presents the experimental results. First, 881 in the image-level category, we find that utilizing TCN can 882 improve the performance over the frame level method, which 883 demonstrates the effectiveness of TCN in modelling temporal 884 dependency. Second, the person-level methods perform better 885 than the whole frame based methods. This is because the later 886 ones can better focus on the action performer, which provides 887 more discriminative power of action. Finally, we observe that 888 adopting our semantic-preserving attention and GCN model 889 can further improve the performance, which indicates the 890 discriminative power of features learned by our proposed 891 method. We also show several action segmentation results in 892 supplementary material for visualization. 893

Confusion Matrices on the CA dataset



class of Walking and Crossing as the same class of Moving as suggested in [6]. (a) Baseline - HDTM. (b)  $Ours^{\dagger} + GCN_{-SPA+KD}$ . (c) Baseline - HDTM. (d)  $Ours^{\dagger} + GCN_{-SPA+KD}$ .

# TABLE V



Method	Accuracy
STL <sup>‡</sup> [60]	77.4
MTCAR <sup>‡ ?</sup> [7]	83.0
Ours-teacher*	79.3
Ours-teacher* + GCN	79.8
Ours-teacher	70.2
baseline-HDTM [2]	57.0
Ours _ SA	57.3
Ours _ SPA	58.3
Ours $-SPA + KD$	58.5
Ours <sup>†</sup> – SA	76.2
Ours <sup>†</sup> – SPA	77.3
Ours <sup>†</sup> – SPA + KD	77.5
Ours + GCN $_{-SA}$	57.9
Ours + GCN _ SPA+KD	58.6
Ours <sup>†</sup> + GCN $_{-SA}$	76.8
Ours <sup>†</sup> + GCN $_{-SPA}$ + KD	78.1

# TABLE VI

COMPARISON OF THE GROUP ACTIVITY SEGMENTATION ACCURACY (%) ON THE UNTRIMMED VOLLEYBALL DATASET

Method	Category	F1 score
VGG16 [64]	Image level	41.74
TCN [12]	Image level	45.17
TCN_SA	Person level	56.06
TCN_SPA+KD	Person level	57.59
TCN-GCN_SPA+KD	Person level	59.49

## H. Analysis on the Influence of Caption Quality 894

Captions, which are a sets of individual words of actions in 895 this paper, are utilized during three stages in our method: 896

Stage 1: Finetuning the DCNN and LSTM network, and 897 extracting the features of individual actions. 898

Stage 2: Training the Teacher network. 899

Stage 3: Guiding the training process of the Student net-900 work. 901

The Stage 1 is a common process in most deep-learning 902 based methods [2], [3], [6] and the Stage 2 is an intermediate 903 process of our method. The Stage 3 is what we should pay 904

Confusion Matrices on the CAE dataset



Fig. 9. Comparison of Confusion Matrices on CA [60] and CAE dataset [8]. † is defined in the caption of Table I. For the CA datset, we merge the

TABLE VII ANALYSIS ON THE INFLUENCE OF INFERIOR CAPTIONS ON THE SPLIT2 OF CHOI'S DATASET

Method	Accuracy (%)	Influence (%)
Teacher*	79.2	-
Teacher*-new	58.5	-20.7 (Stage 2)
Student	74.4	-
Student-new	60.8	-13.6 (Stage 1)
Student-new_ $SPA + KD$	59.1	-1.7 (Stage 3)

more attention to, as it is the core step of our method and 905 directly influences the final recognition result. 906

In order to further analyze the influence of the caption 907 quality, we conducted the experiments on the split2 of Choi's 908 dataset. We randomly selected 50% captions in the training 909 sets and assigned random single action labels to them. In this 910 way, the caption quality will become inferior. 911

Table VII presents the comparison between results on 912 the original setting (Teacher\*, Student) and the new 913 setting (Teacher\*-new, Student-new, Student-new $_{-SPA + KD}$ ). 914 We observe that the captions will heavily influence Stage 1 and 915 Stage 2 (The accuracy drop from 74.4% (Student) to 60.8% 916 (Student-new) because the extracted features became inferior). 917 In comparison, the decrease caused by our method (Stage 3) is 918 slight, which shows its robustness to the low quality captions. 919 The intuition of our method's robustness lies in two folds. 920 First, as the Teacher Network is trained with noisy input 921 labels, the semantics-preserving attention would tend to learn 922 to deal with such noise. Second, knowledge distillation from 923 Teacher Network provides additional soft labels for training 924 Student Network, which will inevitably cause the decrease of 925 the Student Network if the Teacher Network is noisy. But 926 with ground-truth group activity label as direct supervision, 927 this decrease in performance is relieved and won't hurt the 928 final result too much. 929

# I. Analysis on the Computational Time

There are some real-world applications for group activity 931 recognition, e.g., sports video analysis and traffic surveil-932 lance, which require recognizing the activity in real time. 933 Therefore, we are motivated to investigate the time cost 934 of our approach. Table VIII shows the computational time 935

Computational Time Analysis on the Volleyball Dataset.  $^\dagger$  Is Defined in the Caption of Table I

Training Process (Based on Dataset)	Time (h)
Train Teacher Network	0.36
Train DCNN and LSTM for RGB Images	11.50
Extract Features for RGB Images	0.46
Train GCN, Attention Module and BLSTM	1.00
Compute Optical Flow	61.48
Train DCNN and LSTM for Optical Flow	11.50
Extract Features(OF)	0.46
<sup>†</sup> Train GCN, Attention Module and BLSTM	1.16
Testing Process (Based on Single Frame)	Time (ms)
Extract Features for RGB Images	$8.01 \times 12$ (people)
Activity Recognition (10 Frames)	13.93
Compute Optical Flow	434.65
Extract Features for Optical Flow	$8.01 \times 12$ (people)
<sup>†</sup> Activity Recognition (10 Frames)	26.45

# TABLE IX

Comparison of the Computational Time (s) of Different Methods on the Volleyball Dataset. The Results Are Based on a Clip With 10 Frames.<sup>†</sup> Denotes That the Results Are Based on the Inputs With RGB Images and Optical Flows

SBGAR [4]	HDTM [2]	Ours_SPA + KD	Ours+GCN <sub><math>-SPA + KD</math></sub>
-	0.950	0.968	0.983
$1.0966^{\dagger}$	6.207†	6.227†	6.295 <sup>†</sup>

analysis of our method. The training data were based on one
run while the testing data were averaged over five runs on
the Volleyball dataset. We did not include the time to detect
individual players as we utilized the off-the-shelf tracklets
provided by [2].

Without utilizing optical flow, it required about 0.36 + 11.50 + 0.46 + 1.00 = 13.32h to train the SPTS + GCN. For a video clip with 10 frames, it took  $10 \times (8.01 \times 12) + 13.93 =$ 983.14*ms*(0.983*sec*) to predict the group activity label. Moreover, training the Teacher Network was about 0.36 h, only 946 2.70% of the entire training time.

When combining the optical flow, the training phase lasted about  $0.36 + 61.48 + 2 \times (11.5 + 0.46) + 1.16 = 86.92h$ while predicting the label of a video clip took  $10 \times (434.65 + 8.01 \times 12 \times 2) + 26.45 = 6295.35ms(6.295sec)$ . The reason why combining the optical flow is relatively slow is that, we employed the Flownet 2.0 model with the best performance and highest computational time cost in [58].

Table IX presents the computational time comparison with 954 state-of-the-arts. The result of SBGAR is reported from [4], 955 and the others are based on our implementation. On one hand, 956 we find that when combining optical flow, the SBGAR is more 957 efficient and the reason are two folds. (1) The optical flow 958 computation time of SBGAR on a single image is much faster 959 than ours (0.022s vs 0.435s) due to the difference between 960 the methods for calculating optical flow. (2) SBGAR directly 961 takes the whole frames as inputs while our method is based 962 on the a set of tracklets. On the other hand, compared with 963 the baseline approach HDTM [2], the increased time cost 964 of  $Ours_{-SPA + KD}$  and  $Ours_{-SPA + KD}$  are slight, which 965 illustrates the efficiency of our methods. 966

# V. FUTURE WORKS

There are some interesting directions for future works:

- 1) Designing different formulations of GCN for group 969 activity recognition. For example, one is to use a single 970 graph with temporal information. Concretely, we can 971 first perform temporal pooling (e.g., max-pooling or 972 attention-pooling) over the features of individual person 973 and adjacency matrices of different frames, and then 974 construct a single graph and feed it into the GCN model. 975 Another one, which is inspired by [47], is to build a 976 spatial-temporal graph. In this way, features of different 977 people in different frames will be organized in a unified 978 graph, and the final bidirectional LSTM layer in our 979 model can be removed. However, as the scale of the 980 spatial-temporal graph is much larger, other efforts on 981 efficient modeling need to be devoted. 982
- 2) Transferring knowledge in the graph between the Student and Teacher network.<sup>2</sup>
- Employing our method for the tasks like image/video caption or visual question answering (VQA), which lie in the interaction area of the natural language domain and computer vision domain.
- 4) Exploring different variants in [58] and other optical flow estimation algorithms to achieve a better trade-off between the accuracy and efficiency.
   990 991 990 991

# VI. CONCLUSIONS

In this paper, we have presented a Semantics-Preserving 993 Teacher-Student (SPTS) architecture for group activity recog-994 nition in videos. The proposed method has explored the 995 attention knowledge in the semantic domain and employed 996 it to guide the learning process in appearance domain, 997 which explicitly exploits the attention information of the 998 group people. Moreover, we have strengthened our SPTS 999 by incorporating with two graph convolutional modules to 1000 reason the relationship among different people. Furthermore, 1001 we have extended our approach on action segmentation task 1002 for untrimmed videos and demonstrated its effectiveness. 1003 Extensive experimental results on four datasets have shown the 1004 superior performance of our proposed method in comparison 1005 with the state-of-the-arts. 1006

# Acknowledgement

1007

1011

The authors would like to thank Peiyang Li, 1008 Danyang Zhang, Yu Zheng, Simin Wang, Yongming Rao, and 1009 Tianmin Shu for their generous help. 1010

# References

- H. Hajimirsadeghi, W. Yan, A. Vahdat, and G. Mori, "Visual recognition by counting instances: A multi-instance cardinality potential kernel," in *Proc. CVPR*, Jun. 2015, pp. 2596–2605.
- M. S. Ibrahim, S. Muralidharan, Z. Deng, A. Vahdat, and G. Mori, "A hierarchical deep temporal model for group activity recognition," in *Proc. CVPR*, Jun. 2016, pp. 1971–1980.
   1017

<sup>2</sup>We have made some attempts on this direction, see supplementary material for details.

967

968

983

984

[3] T. Shu, S. Todorovic, and S.-C. Zhu, "CERN: Confidence-energy recurrent network for group activity recognition," in Proc. CVPR, Jul. 2017, pp. 4255-4263.

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1032

1051

- X. Li and M. C. Chuah, "SBGAR: Semantics based group activity [41 recognition," in Proc. ICCV, Oct. 2017, pp. 2895-2904
- T. M. Bagautdinov, A. Alahi, F. Fleuret, P. Fua, and S. Savarese, "Social scene understanding: End-to-end multi-person action localization and collective activity recognition," in Proc. CVPR, Jul. 2017, pp. 3425-3434
- M. Wang, B. Ni, and X. Yang, "Recurrent modeling of interaction [6] context for collective activity recognition," in Proc. CVPR, Jul. 2017, pp. 7408-7416.
- 1030 W. Choi and S. Savarese, "A unified framework for multi-target tracking and collective activity recognition," in Proc. ECCV, 2012, pp. 215-230. 1031
- W. Choi, K. Shahid, and S. Savarese, "Learning context for collective activity recognition," in Proc. CVPR, Jun. 2011, pp. 3273-3280. 1033
- T. Lan, Y. Wang, W. Yang, S. N. Robinovitch, and G. Mori, "Dis-1034 [9] criminative latent models for recognizing contextual group activities," 1035 IEEE Trans. Pattern Anal. Mach. Intell., vol. 34, no. 8, pp. 1549-1562, 1036 Aug. 2012. 1037
- [10] T. Shu, D. Xie, B. Rothrock, S. Todorovic, and S.-C. Zhu, "Joint 1038 inference of groups, events and human roles in aerial videos," in Proc. 1039 1040 CVPR, Jun. 2015, pp. 4576-4584.
- Y. Tang, P. Zhang, J.-F. Hu, and W.-S. Zheng, "Latent embeddings for 1041 [11] 1042 collective activity recognition," in Proc. AVSS, Aug./Sep. 2017, pp. 1-6.
- C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager, "Temporal [12] 1043 1044 convolutional networks for action segmentation and detection," in Proc. 1045 CVPR, Jul. 2017, pp. 1003-1012.
- [13] Y. Tang, Z. Wang, P. Li, J. Lu, M. Yang, and J. Zhou, "Mining semantics-1046 preserving attention for group activity recognition," in Proc. 26th ACM 1047 Int. Conf. Multimedia, 2018, pp. 1283-1291. 1048
- H. Wang, H. Kläser, A. Schmid, and C.-L. Liu, "Action recognition by [14] 1049 dense trajectories," in Proc. CVPR, Jun. 2011, pp. 3169-3176. 1050
- W. Du, Y. Wang, and Y. Qiao, "Recurrent spatial-temporal attention [15] network for action recognition in videos," IEEE Trans. Image Process., 1052 vol. 27, no. 3, pp. 1347-1360, Mar. 2018.
- S. Ji, W. Xu, M. Yang, and K. Yu, "3D convolutional neural networks 1054 [16] 1055 for human action recognition," IEEE Trans. Pattern Anal. Mach. Intell., vol. 35, no. 1, pp. 221-231, Jan. 2013. 1056
- W. Hu, B. Wu, P. Wang, C. Yuan, Y. Li, and S. J. Maybank, "Context-1057 [17] dependent random walk graph kernels and tree pattern graph matching 1058 kernels with applications to action recognition," IEEE Trans. Image 1059 Process., vol. 27, no. 10, pp. 5060-5075, Oct. 2018. 1060
- [18] K. Simonyan and A. Zisserman, "Two-stream convolutional networks 1061 for action recognition in videos," in Proc. NIPS, 2014, pp. 568-576. 1062
- X. Chang, W.-S. Zheng, and J. Zhang, "Learning person-person inter-action in collective activity recognition," *IEEE Trans. Image Process.*, 1063 [19] 1064 vol. 24, no. 6, pp. 1905–1918, Jun. 2015. 1065
- [20] M. S. Ibrahim and G. Mori, "Hierarchical relational networks for group 1066 activity recognition and retrieval," in Proc. ECCV, 2018, pp. 721-736. 1067
- [21] J. K. Tsotsos, S. M. Culhane, W. Y. K. Wai, Y. Lai, N. Davis, and 1068 F. Nuflo, "Modeling visual attention via selective tuning," Artif. Intell., 1069 vol. 78, nos. 1-2, pp. 507-545, 1995. 1070
- [22] R. A. Rensink, "The dynamic representation of scenes," Vis. Cognition, 1071 vol. 7, nos. 1-3, pp. 17-42, Oct. 2010. 1072
- D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by 1073 [23] jointly learning to align and translate," in Proc. ICLR, 2014, pp. 1-15. 1074
- [24] A. Vaswani et al., "Attention is all you need," in Proc. NIPS, 2017, 1075 pp. 6000-6010. 1076
- J. Yang et al., "Neural aggregation network for video face recognition," 1077 [25] in Proc. CVPR, Jul. 2017, pp. 5216-5225. 1078
- 1079 [26] Y. Rao, J. Lu, and J. Zhou, "Attention-aware deep reinforcement learning for video face recognition," in Proc. ICCV, Oct. 2017, pp. 3951-3960. 1080
- A. Haque, A. Alahi, and L. Fei-Fei, "Recurrent attention models 1081 [27] for depth-based person identification," in Proc. CVPR, Jun. 2016, 1082 pp. 1229-1238. 1083
- [28] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Is object localization for 1084 1085 free?—Weakly-supervised learning with convolutional neural networks." 1086 in Proc. CVPR, Jun. 2015, pp. 685-694.
- S. Zagoruyko and N. Komodakis, "Paying more attention to attention: 1087 [29] 1088 Improving the performance of convolutional neural networks via attention transfer," in Proc. ICLR, 2017, PP. 1-13. 1089
- Z. Guo, L. Gao, J. Song, X. Xu, J. Shao, and H. T. Shen, "Attention-1090 [30] based LSTM with semantic consistency for videos captioning," in *Proc.* 1091 24th ACM Int. Conf. Multimedia, 2016, pp. 357-361. 1092

- [31] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention 1093 networks for image question answering," in Proc. CVPR, Jun. 2016, 1094 pp. 21-29. 1095
- [32] J. Liu, G. Wang, L.-Y. Duan, K. Abdiyeva, and A. C. Kot, "Skeleton-1096 based human action recognition with global context-aware atten-1097 tion LSTM networks," IEEE Trans. Image Process., vol. 27, no. 4, 1098 pp. 1586-1599, Apr. 2018, 1099
- [33] S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu, "An end-to-end spatio-1100 temporal attention model for human action recognition from skeleton 1101 data," in Proc. AAAI, 2017, pp. 4263-4270. 1102
- [34] G. E. Hinton, O. Vinvals and J. Dean, "Distilling the knowledge in a 1103 neural network," in Proc. NIPSW, 2014. 1104
- A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and [35] 1105 Y. Bengio, "FitNets: Hints for thin deep nets," in Proc. ICLR, 2014, 1106 pp. 1–13. 1107
- [36] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: 1108 Fast optimization, network minimization and transfer learning," in Proc. 1109 CVPR, Jul. 2017, pp. 7130-7138. 1110
- T. Chen, I. J. Goodfellow, and J. Shlens, "Net2Net: Accelerating learning [37] 1111 via knowledge transfer," in Proc. ICLR, 2015, pp. 1-12. 1112
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, [38] 1113 "The graph neural network model," IEEE Trans. Neural Netw., vol. 20, 1114 no. 1, pp. 61-80, Jan. 2009. 1115
- [39] T. N. Kipf and M. Welling, "Semi-supervised classification with graph 1116 convolutional networks," in Proc. ICLR, 2017, pp. 1-14. 1117
- [401 M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural 1118 networks on graphs with fast localized spectral filtering," in Proc. NIPS, 1119 2016, pp. 3844-3852. 1120
- [41] J. Bruna, W. Zaremba, A. Szlam, and Y. Lecun, "Spectral networks and 1121 locally connected networks on graphs," in Proc. ICLR, 2014, pp. 1-14. 1122
- D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: [42] 1123 1124 Extending high-dimensional data analysis to networks and other irregular 1125 domains," IEEE Signal Process. Mag., vol. 30, no. 3, pp. 83-98, 1126 May 2013. 1127
- X. Qi, R. Liao, J. Jia, S. Fidler, and R. Urtasun, "3D graph neural [43] 1128 networks for RGBD semantic segmentation," in Proc. ICCV, Oct. 2017, 1129 pp. 5209-5218. 1130
- X. Chen, L. Li, L. Fei-Fei, and A. Gupta, "Iterative visual reasoning [44] 1131 beyond convolutions," in Proc. CVPR, Jun. 2018, pp. 7239-7248. 1132
- [45] X. Wang, Y. Ye, and A. Gupta, "Zero-shot recognition via seman-1133 tic embeddings and knowledge graphs," in Proc. CVPR, Jun. 2018, 1134 pp. 6857-6866. 1135
- Z. Wang, T. Chen, J. Ren, W. Yu, H. Cheng, and L. Lin, "Deep reasoning [46] 1136 with knowledge graph for social relationship understanding," in Proc. 1137 IJCAI, 2018, pp. 1021-1028. 1138
- [47] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional 1139 networks for skeleton-based action recognition," in Proc. AAAI, 2018, 1140 pp. 7444-7452. 1141
- [48] Y. Tang, Y. Tian, J. Lu, P. Li, and J. Zhou, "Deep progressive reinforce-1142 ment learning for skeleton-based action recognition," in Proc. CVPR, 1143 Jun. 2018, pp. 5323-5332 1144
- [49] C. Li, Z. Cui, W. Zheng, C. Xu, and J. Yang, "Spatio-temporal graph 1145 convolution for skeleton based action recognition," in Proc. AAAI, 2018, 1146 pp. 3482-3489. 1147
- [50] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classifica-1148 tion with deep convolutional neural networks," in Proc. NIPS, 2012, 1149 pp. 84-90. 1150
- [51] P. J. Werbos, "Backpropagation through time: What it does and how to 1151 do it," Proc. IEEE, vol. 78, no. 10, pp. 1550-1560, Oct. 1990. 1152
- D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: [52] 1153 The missing ingredient for fast stylization," CoRR, abs/1607.08022, 1154 Jul. 2016. 1155
- [53] X. Xiang, Y. Tian, A. Reiter, G. D. Hager, and T. D. Tran, "S3D: 1156 Stacking segmental P3D for action quality assessment," in Proc. ICIP, 1157 Oct. 2018, pp. 928-932. 1158
- [54] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, X. Tang, and D. Lin, "Temporal 1159 action detection with structured segment networks," in Proc. ICCV, 1160 Oct. 2017, pp. 2933-2942. 1161
- H. Xu, A. Das, and K. Saenko, "R-C3D: Region convolutional 3D [55] 1162 network for temporal activity detection," in Proc. ICCV, Oct. 2017, 1163 pp. 5794-5803. 1164
- Y. Tang et al., "COIN: A large-scale dataset for comprehensive instruc-[56] 1165 tional video analysis," in Proc. CVPR, 201, pp. 1-10. 1166

- [57] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards realtime object detection with region proposal networks," in *Proc. NIPS*, 2015, pp. 91–99.
- [58] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox,
  "FlowNet 2.0: Evolution of optical flow estimation with deep networks,"
  in *Proc. CVPR*, Jul. 2017, pp. 1647–1655.
- [59] M. S. Ibrahim, S. Muralidharan, Z. Deng, A. Vahdat, and G. Mori,
   "Hierarchical deep temporal models for group activity recognition,"
   *CoRR*, abs/1607.02643, Jul. 2016.
- [60] W. Choi, K. Shahid, and S. Savarese, "What are they doing?: Collective activity classification using spatio-temporal relationship among people,"
   in *Proc. ICCVW*, Sep./Oct. 2009, pp. 1282–1289.
- [61] Z. Deng, A. Vahdat, H. Hu, and G. Mori, "Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition," in *Proc. CVPR*, Jun. 2016, pp. 4772–4781.
- [62] A. Fathi, X. Ren, and J. M. Rehg, "Learning to recognize objects in egocentric activities," in *Proc. CVPR*, Jun. 2011, pp. 3281–3288.
- [63] S. Stein and S. J. McKenna, "Combining embedded accelerometers with computer vision for recognizing food preparation activities," in *Proc. ACM Int. Joint Conf. Pervasive Ubiquitous Comput.*, 2013, pp. 729–738.
- [64] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. ICLR*, 2015, pp. 1–14.
- <sup>1190</sup> [65] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [66] S. Biswas and J. Gall, "Structural recurrent neural network (SRNN) for
   group activity analysis," in *Proc. WACV*, Mar. 2018, pp. 1625–1632.

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211



Yansong Tang received the B.S. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2015, where he is currently pursuing the Ph.D. degree with the Department of Automation. His research lies in computer vision, especially multi-modal action recognition and egocentric vision analytics.

Jiwen Lu (M'11–SM'15) received the B.Eng. degree in mechanical engineering and the M.Eng. degree in electrical engineering from the Xi'an University of Technology, Xi'an, China, in 2003 and 2006, respectively, and the Ph.D. degree in electrical engineering from Nanyang Technological University, Singapore, in 2012. He is currently an Associate Professor with the Department of Automation, Tsinghua University, Beijing, China. His current research interests include computer vision, pattern recognition, and machine learning. He has

authored/coauthored over 200 scientific papers in these areas, where over 60 of 1212 1213 them are the IEEE TRANSACTIONS papers (including 13 T-PAMI papers) and 50 of them are CVPR/ICCV/ECCV/NIPS papers. He is a member of the Mul-1214 timedia Signal Processing Technical Committee and the Information Forensics 1215 and Security Technical Committee of the IEEE Signal Processing Society, and 1216 a member of the Multimedia Systems and Applications Technical Committee 1217 and the Visual Signal Processing and Communications Technical Committee 1218 of the IEEE Circuits and Systems Society. He was a recipient of the National 1219 1000 Young Talents Program of China in 2015, and the National Science 1220 Fund of China for Excellent Young Scholars in 2018, respectively. He serves 1221 as the Co-Editor-of-Chief for the Pattern Recognition Letters, an Associate 1222 Editor of the IEEE TRANSACTIONS ON IMAGE PROCESSING, the IEEE 1223 TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY, 1224 the IEEE TRANSACTIONS ON BIOMETRICS, BEHAVIOR, AND IDENTITY 1225 1226 SCIENCE, and Pattern Recognition.



Zian Wang is currently pursuing the B.S. degree 1227 with the Department of Automation, Tsinghua University, Beijing, China. His research interests include 1229 computer vision and machine learning. 1230



Ming Yang (M'08) received the B.E. and M.E. 1231 degrees in electronic engineering from Tsinghua 1232 University, Beijing, China, in 2001 and 2004, respec-1233 tively, and the Ph.D. degree in electrical and com-1234 puter engineering from Northwestern University, 1235 Evanston, IL, USA, in 2008. From 2004 to 2008, 1236 he was a Research Assistant with the Computer 1237 Vision Group, Northwestern University. After his 1238 graduation, he joined NEC Laboratories America, 1239 Cupertino, CA, USA, where he was a Senior 1240 Researcher. He was a Research Scientist in AI 1241

Research at Facebook (FAIR) from 2013 to 2015. He is currently the Co-Founder and VP of software at Horizon Robotics, Inc. His research interests include computer vision, machine learning, face recognition, large scale image retrieval, and intelligent multimedia content analysis. He is the author of over 50 peer-reviewed publications in prestigious international journals and conferences, which have been cited over 9400 times.



Jie Zhou (M'01-SM'04) received the B.S. and 1248 M.S. degrees from the Department of Mathemat-1249 ics, Nankai University, Tianjin, China, in 1990 1250 and 1992, respectively, and the Ph.D. degree from 1251 the Institute of Pattern Recognition and Artificial 1252 Intelligence, Huazhong University of Science and 1253 Technology (HUST), Wuhan, China, in 1995. Since 1254 then, he has served as a Post-Doctoral Fellow at 1255 the Department of Automation, Tsinghua University, 1256 Beijing, China, until 1997. Since 2003, he has been 1257 a Full Professor with the Department of Automation, 1258

Tsinghua University. His research interests include computer vision, pattern 1259 recognition, and image processing. In recent years, he has authored over 1260 100 papers in peer-reviewed journals and conferences. Among them, over 1261 30 papers have been published in top journals and conferences, such as 1262 the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTEL-1263 LIGENCE, the IEEE TRANSACTIONS ON IMAGE PROCESSING, and CVPR. 1264 He is an Associate Editor of the IEEE TRANSACTIONS ON PATTERN ANALY-1265 SIS AND MACHINE INTELLIGENCE and two other journals. He received the 1266 National Outstanding Youth Foundation of China Award. 1267