# A Fully Connected Model for Consistent Collective Activity Recognition in Videos

Takuhiro Kaneko\*, Masamichi Shimosaka, Shigeyuki Odashima, Rui Fukui, Tomomasa Sato

Department of Mechano-Informatics, the University of Tokyo, Tokyo, Japan

#### Abstract

We propose a novel method for consistent collective activity recognition in video images. Collective activities are activities performed by multiple persons, such as queuing in a line, talking together, and waiting at an intersection. Since it is often difficult to differentiate between these activities using the appearance of only an individual person, the models proposed in recent studies exploit the contextual information of other people nearby. However, these models do not sufficiently consider the spatial and temporal consistency in a group (e.g., they consider the consistency in only the adjacent area), and therefore, they cannot effectively deal with temporary misclassification or simultaneously consider multiple collective activities in a scene. To overcome this drawback, this paper describes a method to integrate the individual recognition results via fully connected conditional random fields (CRFs), which consider all the interactions among the people in a video clip and alter the interaction strength in accordance with the degree of their similarity. Unlike previous methods that restrict the interactions among the

Preprint submitted to Pattern Recognition Letters

October 23, 2014

<sup>\*</sup>Corresponding author

Email address: kaneko@ics.t.u-tokyo.ac.jp (Takuhiro Kaneko)

people heuristically (e.g., within a constant area), our method describes the "multi-scale" interactions in various features, i.e., position, size, motion, and time sequence, in order to allow various types, sizes, and shapes of groups to be treated. Experimental results on two challenging video datasets indicate that our model outperforms not only other graph topologies but also state-of-the art models.

Keywords:

Collective activity recognition, Fully connected model, CRFs, Spatial and temporal consistency

# 1 1. Introduction

Vision-based human activity recognition is of scientific and practical im-2 portance, and has been actively studied in the research field of computer 3 vision. Many previous studies focused on recognizing actions performed by a single person in a video clip (Blank et al., 2005; Niebles et al., 2006; Schuldt 5 et al., 2004). However, in real-world applications, such as surveillance moni-6 toring, the previous methods are inapplicable, since human actions are rarely 7 performed by a single person, but instead by multiple persons. For exam-8 ple, it is difficult to differentiate between the activities of the two persons 9 shown in Fig. 1(a), by considering the appearance of the individual person. 10 In order to recognize activities performed by multiple persons, which we call 11 "collective activities," it is necessary to exploit the contextual information of 12 the people nearby. When we have identified the activities of people nearby, 13 it immediately becomes clear that the left person in Fig. 1(a) is queuing and 14 the right person is talking, as shown in Fig. 1(b). 15



Figure 1: Useful contexts for collective activity recognition. It is often difficult to differentiate between collective activities by the appearance of only an individual person (a). When we have identified the activities of people nearby, it immediately becomes clear that the left person is queuing and the right person is talking (b).

In some recent studies, methods have been proposed for collective ac-16 tivity recognition using the contextual information of people nearby. Choi 17 et al. (2009), Lan et al. (2010a), and Kaneko et al. (2012b) encoded the 18 contextual information by exploiting the feature descriptors extracted from 19 a focal person and his/her surrounding area. These descriptors are more ef-20 fective than feature descriptors without contexts (e.g., histogram of oriented 21 gradients (HOG) (Dalal and Triggs, 2005)). However, in the models, the 22 activity of each person is classified independently, and therefore, the spatial 23 and temporal consistency in a group is not always ensured. 24

In order to obtain this consistency, the question "Which people are in the same group?" must be answered, and an activity in each group must be

optimized. To answer the question, Amer and Todorovic (2011) optimized ac-27 tivities around deformable grids, while Lan et al. (2010b), Choi et al. (2011), 28 Choi and Savarese (2012), Khamis et al. (2012a,b) used graph structures 29 that describe the interactions between persons. However, the models used 30 in these studies cannot describe the "multi-scale" interactions in various fea-31 tures, such as position, size, motion and time sequence, although there exist 32 various types, sizes, and shapes of groups, as shown in Fig. 2. The model pro-33 posed by Amer and Todorovic (2011) depended on the density and position 34 of the grids, and therefore, it was difficult to exploit long-range relationships. 35 In the model proposed by Lan et al. (2010b), the person-person interactions 36 were latent and learned automatically; however, their model was restricted 37 to modeling contextual information in a single frame, and was not designed 38 such that temporal consistency was ensured. Choi et al. (2011) and Khamis 39 et al. (2012a,b) considered temporal consistency for a person or group; how-40 ever, in their model, the person-person interactions that were considered were 41 restricted only in consecutive frames to compute reasonably. The results of 42 these models are likely to be affected by temporary misclassification. Choi 43 and Savarese (2012) exploited a hierarchical model to classify collective activ-44 ities jointly; however, they assumed there exists only one collective activity 45 in a certain time frame. Therefore, the method cannot model multiple collec-46 tive activities in a scene, such as that shown in Fig. 2(b), where some persons 47 are waiting at a street intersection, while others are crossing. Considering 48 real-world applications, such as surveillance monitoring, this assumption is 49 not natural. 50



In contrast, our proposed model describes the "multi-scale" interactions



Figure 2: Which people are in the same group? For dividing people into groups, various criteria, such as (a) position, (b) size, (c) motion, and (d) time sequence, can be used.

in various features, i.e., position, size, motion, and time sequence. This 52 means that our model is able not only to describe the long-range relationships 53 among people in both time and space, but also to consider multiple collective 54 activities in a certain time frame. In particular, we use fully connected 55 conditional random fields (CRFs), which consider all the interactions among 56 the people in a video clip, and alter the interaction strength according to the 57 degree of their similarity. This model is able to represent the various features 58 over a "multi-scale" in a single unified model. In general, the calculation cost 59 of a fully connected model is intractable when strict estimation is conducted; 60 however, the cost is reduced to linear in the number of detected persons using 61 a highly efficient approximation method in which the pairwise potentials are 62 modeled using Gaussian kernels (Krähenbühl and Koltun, 2011). 63

We summarize the main contributions of this paper. (1) We propose a 64 novel method for consistent collective activity recognition in video images 65 using a fully connected model. In the model, we do not restrict the person-66 person interactions that are considered heuristically, but instead consider all 67 the interactions among the people in a video clip. (2) We describe the person-68 person interactions over the multi-scale, using various features: position, 69 size, motion, and time sequence. The interaction strength among the people 70 is altered according to the degree of their similarity in the features. (3) 71 We perform the inference with linear complexity in the number of detected 72 persons, using an approximation method in which the pairwise potential are 73 modeled using Gaussian kernels. (4) We evaluate our model on two publicly 74 available datasets. The experimental results show that our fully connected 75 model outperforms other graph structures, such as the unary only model, 76 and the adjacently connected model, as well as state-of-the art models (Choi 77 et al., 2009, 2011; Khamis et al., 2012a,b; Lan et al., 2010a; Kaneko et al., 78 2012b). Portions of this paper appeared previously in Kaneko et al. (2012a). 79 In this paper, we additionally evaluate our model on the dataset (Choi and 80 Savarese, 2012) and report the results to make our contribution stronger. 81 We also present comparisons with other graph structures, and an additional 82 analysis of the qualitative results, having clarified the characteristic of our 83 model. Moreover, we evaluate a novel combination of our fully connected 84 model and a state-of-the art feature descriptor (Kaneko et al., 2012b), and 85 report the results of a comparison of our proposed model and state-of-the art 86 models. 87

88

The rest of the paper is organized as follows. First, in Section 2, we

present our framework of consistent collective activity recognition in video images. Next, in Section 3, the details of learning and inference of the model are given. In Section 4, we report our experimental results quantitatively and qualitatively. Finally, we summarize our paper and present our conclusions in Section 5.

# <sup>94</sup> 2. Consistent Collective Activity Recognition with Fully Connected <sup>95</sup> CRFs

## 96 2.1. Model Overview

The main goal of our study is to ensure the spatial and temporal consis-97 tency of the activity in each group in collective activity recognition. For this 98 purpose, our method uses CRFs (Lafferty et al., 2001). CRFs are a proba-99 bilistic framework for labeling and segmenting structured data and able to 100 deal with arbitrary dependencies on the observation sequence in a single uni-101 fied model (He et al., 2004; Shotton et al., 2006). Specifically, in order to 102 handle "multi-scale" interactions, our method uses fully connected CRFs. In-103 stead of specifying the interactions among the people heuristically, our model 104 describes the interactions in position, size, motion, and time sequence as the 105 variable potentials, according to the degree of their similarity, in order to 106 allow various types, sizes, and shapes of groups to be treated. 107

We now give a brief overview of our model. In the preprocessing step, persons in a video clip are found. Next, features (e.g., HOG (Dalal and Triggs, 2005) and optical flow) are extracted from the detected bounding box. The unary potentials and the pairwise potentials are calculated using the features, and integrated via fully connected CRFs. We present the technical <sup>113</sup> details of our model in the following sections.

#### 114 2.2. Model Formulation

Fully Connected CRFs Model. Given a video clip, our method first detects 115 persons using an efficient human detector. In our implementation, we use the 116 approach of Felzenszwalb et al. (2008). The observed data of the detected 117 persons are defined as  $\boldsymbol{x} = \{x_1, ..., x_N\}$ , where  $x_i$  is the observed data of the 118 *i*-th person and N is the number of detected persons in the video clip. Let 119 the corresponding activity classes be given by  $\boldsymbol{y} = \{y_1, ..., y_N\}$ . The domain 120 of each variable  $y_i$  is a set of activity classes  $\mathcal{L} = \{l_1, ..., l_K\}$ , where K is the 121 number of activity classes. A conditional random field (x, y) is characterized 122 by a Gibbs distribution: 123

$$P(\boldsymbol{y}|\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{x})} \exp(-E(\boldsymbol{y}|\boldsymbol{x})), \qquad (1)$$

where  $Z(\boldsymbol{x}) = \sum_{\boldsymbol{y}'} \exp(-E(\boldsymbol{y}'|\boldsymbol{x}))$  is the partition function that normalizes the distribution, and  $E(\boldsymbol{y}|\boldsymbol{x})$  is the Gibbs energy, which is associated with a configuration  $\boldsymbol{y}$  conditioned on  $\boldsymbol{x}$ . In the fully connected pairwise CRF model, the Gibbs energy is defined as

$$E(\boldsymbol{y}|\boldsymbol{x}) = \sum_{i} \underbrace{\psi_{u}(y_{i})}_{\text{unary potential}} + \sum_{i} \sum_{j>i} \underbrace{\psi_{p}(y_{i}, y_{j})}_{\text{pairwise potential}}, \quad (2)$$

where  $\psi_u(y_i)$  is the unary potential and  $\psi_p(y_i, y_j)$  is the pairwise potential. For notational convenience, we omit the conditioning in the rest of this paper, and use  $\psi_c(\boldsymbol{y}_c)$  to denote  $\phi_c(\boldsymbol{y}_c|\boldsymbol{x})$  for each clique c.

Unary Potential. The unary potential  $\psi_u(y_i)$  is computed independently for each person by a classifier that produces a distribution over the activity label  $_{133}$   $y_i$  given a contextual feature descriptor

$$\psi_u(y_i) = -\log(P(y_i)),\tag{3}$$

where  $P(y_i)$  represents the probability that the activity of *i*-th person is  $y_i$ .  $P(y_i)$  is calculated by normalizing the classifier scores on the descriptor using a softmax function. In our implementation, we use the descriptors that encode information about not only the action of an individual person, but also the behavior of other people nearby (Choi et al., 2009; Kaneko et al., 2012b). The details are described in Section 2.3.

*Pairwise Potential.* Since the output of the unary classifier for each person is 140 produced independently of the outputs of the classifiers for other people, the 141 recognition result achieved by the unary classifiers alone is generally noisy 142 and inconsistent. To obtain consistency, we exploit the pairwise potential. 143 The pairwise potential  $\psi_p(y_i, y_j)$  represents the interactions between persons. 144 In our fully connected model, the graph structure is the complete graph on 145 y, and the pairwise potential is computed for all the sets of persons in a 146 video clip. The detailed explanation of our graph structure is described in 147 Fig 3. In our model, the pairwise potential is defined as 148

$$\psi_p(y_i, y_j) = \mu(y_i, y_j) k(\boldsymbol{f}_i, \boldsymbol{f}_j), \qquad (4)$$

where  $\mu(y_i, y_j)$  is the label compatibility function given by the Potts model (Boykov and Jolly, 2001):  $\mu(y_i, y_j) = [y_i \neq y_j]$ . It introduces a penalty for similar persons that are assigned different labels. The vectors  $f_i$  and  $f_j$  are feature vectors for the *i*-th and *j*-th persons, and  $k(f_i, f_j)$  is the Gaussian kernel defined by the positions  $p_i$  and  $p_j$ , sizes  $s_i$  and  $s_j$ , motions  $m_i$  and  $m_j$ , 154 times  $t_i$  and  $t_j$ , and weight w:

$$k(\mathbf{f}_{i}, \mathbf{f}_{j}) = w \exp\left(-\frac{|p_{i} - p_{j}|^{2}}{2\theta_{1}^{2}} - \frac{|s_{i} - s_{j}|^{2}}{2\theta_{2}^{2}} - \frac{|m_{i} - m_{j}|^{2}}{2\theta_{3}^{2}} - \frac{|t_{i} - t_{j}|^{2}}{2\theta_{4}^{2}}\right).$$
(5)

The kernel is inspired by the observation that the persons in the same group have similarities in position, size, motion, and time sequence, as illustrated in Fig 2.

It should be noted that we normalize positions and sizes according to 158 the median size of all the persons in the video clip, in order to describe 159 the interaction strength as a relative rather than an absolute quantity. The 160 motion is computed using different methods according to whether a video clip 161 is captured using a moving or a fixed camera. When using a moving camera, 162 the motion is calculated by subtracting the median optical flow without the 163 bounding boxes from the mean optical flow within the bounding box. The 164 former optical flow represents the camera motion, while the latter optical 165 flow represents the person motion in the image. When using a fixed camera, 166 the motion is defined as the mean optical flow within the bounding box. The 167 optical flow is computed using the approach of Sun et al. (2010). 168

#### 169 2.3. Contextual Feature Descriptors

In this section, we describe our method for encoding contextual information into feature descriptors, and calculate the probability in equation (3). We use two descriptors: The *action context (AC) descriptor* proposed by Lan et al. (2010a); and the *combination of the action context and relative action context descriptors (AC-RAC)* that we previously proposed (Kaneko et al., 2012b). In the experiments, we use the AC descriptor as the baseline to evaluate our model in comparison with other graph topologies. In Kaneko et al.
(2012b), it was shown that the AC-RAC descriptor outperforms other previous descriptors (Choi et al., 2009, 2011; Lan et al., 2010a). We integrate our
model with this efficacious descriptor and compare it with state-of-the-art
methods.

The Action Context Descriptor. The AC descriptor (Lan et al., 2010a) is a per-person descriptor; each descriptor is calculated by concatenating the action descriptor, which captures the action of the focal person, and the context descriptor, which captures the behavior of nearby people.

The action descriptor has a bag-of-words style. We employ the person 185 descriptors (e.g., HOG (Dalal and Triggs, 2005)) as the underlying repre-186 sentation. We then train a multiclass SVM classifier associated with action 187 labels. In the experiments, we use a linear SVM implementation of LIBLIN-188 EAR (Fan et al., 2008). Using the score returned by the SVM classifier, the 189 *i*-th person is represented as the K-dimensional vector  $F_i = [S_{1i}, S_{2i}, ..., S_{Ki}]$ , 190 where K is the number of action classes, and  $S_{ki}$  is the score of classifying 191 the *i*-th person to the k-th action class. 192

<sup>193</sup> When the action descriptor has been computed for each person, the con-<sup>194</sup> text descriptor is calculated by integrating the action descriptor of nearby <sup>195</sup> people in the "context region." The context region is further divided into <sup>196</sup> M regions, called "sub-context regions," in space and time, and then the <sup>197</sup> context descriptor is represented as the  $M \times K$  dimensional vector:

$$C_i = [D_{1i}, ..., D_{Mi}]$$

$$= \left[\max_{j \in \mathcal{N}_{1}(i)} S_{1j}, ..., \max_{j \in \mathcal{N}_{1}(i)} S_{Kj}, ..., \max_{j \in \mathcal{N}_{M}(i)} S_{1j}, ..., \max_{j \in \mathcal{N}_{M}(i)} S_{Kj}\right], \quad (6)$$

where  $D_{mi}$  is called the "sub-context descriptor" representing the context in the *m*-th sub-context region of the *i*-th person, and  $\mathcal{N}_m(i)$  denotes the indices of people in the sub-context region.

The AC descriptor of the *i*-th person,  $A_i$ , is computed by concatenating its action descriptor  $F_i$  and its context descriptor  $C_i$ :  $A_i = [F_i, C_i]$ . We then run the multiclass SVM classifier on the AC descriptor associated with activity labels. The classifier scores are normalized using a softmax function, and incorporated as the unary potential in equation (3).

The Relative Action Context Descriptor. The relative action context (RAC) descriptor (Kaneko et al., 2012b) is a refinement of the AC descriptor. Unlike the AC descriptor, the RAC encodes the relative relationship (e.g., if the focal person is facing left and another person is facing right, the relative relationship is defined as facing the opposite direction), and therefore, the descriptor is invariant under a change in the viewpoint, and consistent within the same category of collective activity.

Similarly to Lan et al. (2010b), we define actions by concatenating poses and activities (e.g., talking and facing right). This means that the action descriptor and the sub-context descriptor are  $K(=U \times V)$  dimensional vectors, where U is the number of activity classes and V is the number of pose classes. Using U and V, we redefine the action descriptor  $F_i$ , and the sub-context descriptor  $D_{mi}$  in the AC descriptor:

$$F_{i} = [S_{1i}, S_{2i}, ..., S_{Ki}]$$
  
= [S\_{11i}, S\_{12i}, ..., S\_{uvi}, ..., S\_{UVi}], (7)

$$D_{mi} = \left[\max_{j \in \mathcal{N}_m(i)} S_{1j}, ..., \max_{j \in \mathcal{N}_m(i)} S_{Kj}\right]$$
$$= \left[\max_{j \in \mathcal{N}_m(i)} S_{11j}, \max_{j \in \mathcal{N}_m(i)} S_{12j}, ..., \max_{j \in \mathcal{N}_m(i)} S_{uvj}, ..., \max_{j \in \mathcal{N}_m(i)} S_{UVj}\right].$$
(8)

The RAC descriptor is calculated by shifting the AC descriptor based on the pose of the focal person. First, the pose of the *i*-th person,  $\hat{v}_i$ , is calculated from the person descriptor (e.g., HOG (Dalal and Triggs, 2005)) using a multiclass SVM classifier. In terms of pose  $\hat{v}_i$ , the *i*-th person's relative action descriptor  $\hat{F}_i$  and its relative sub-context descriptor  $\hat{D}_{mi}$  are defined as

$$\hat{F}_{i} = \begin{bmatrix} S_{1\hat{v}_{i}i}, ..., S_{1Vi}, S_{11i}, ..., S_{1(\hat{v}_{i}-1)i}, ..., \\ S_{U\hat{v}_{i}i}, ..., S_{UVi}, S_{U1i}, ..., S_{U(\hat{v}_{i}-1)i} \end{bmatrix},$$

$$\hat{D}_{mi} = \begin{bmatrix} \max_{j \in \mathcal{N}_{m}(i)} S_{1\hat{v}_{i}j}, ..., \max_{j \in \mathcal{N}_{m}(i)} S_{1Vj}, \max_{j \in \mathcal{N}_{m}(i)} S_{11j}, \max_{j \in \mathcal{N}_{m}(i)} S_{1(\hat{v}_{i}-1)j}, ..., \\ \max_{j \in \mathcal{N}_{m}(i)} S_{U\hat{v}_{i}j}, ..., \max_{j \in \mathcal{N}_{m}(i)} S_{UVj}, \max_{j \in \mathcal{N}_{m}(i)} S_{U1j}, \max_{j \in \mathcal{N}_{m}(i)} S_{U(\hat{v}_{i}-1)j} \end{bmatrix}. (10)$$

The relative context descriptor of the *i*-th person,  $\hat{C}_i$ , is computed by concatenating its relative sub-context descriptor:  $\hat{C}_i = [\hat{D}_{1i}, ..., \hat{D}_{Mi}]$ . Finally, the RAC descriptor of the *i*-th person,  $R_i$ , is computed by concatenating its relative action descriptor  $\hat{F}_i$  and its relative context descriptor  $\hat{C}_i$ :  $R_i = [\hat{F}_i, \hat{C}_i]$ .

The Combination of the AC and RAC Descriptors. We now describe the method for combining the AC and RAC descriptors (AC-RAC) (Kaneko et al., 2012b). After extracting the AC and RAC descriptors, we run the multiclass SVM classifier on each of the descriptors associated with activity labels, and transform classifier scores into probabilities via softmax transfor<sup>235</sup> mation. We then combine them via the MAX rule (Hatef et al., 1998):

$$\hat{y}_i = \arg\max_{y_i} P_i(y_i) \text{ s.t. } P_i(y_i) = \max_k P_i(y_i|d_k),$$
 (11)

where  $P_i(y_i)$  is the probability that the activity of the *i*-th person is  $y_i$ , and  $P_i(y_i|d_1)$  and  $P_i(y_i|d_2)$  are the probability calculated from the AC and RAC descriptors, respectively. The probability  $P_i(y_i)$  is incorporated as the unary potential in equation (3).

It should be noted that, when we use the AC-RAC descriptor in our implementation, we employ two post processes (threshold processing and Gaussian filtering) on the AC and RAC descriptors to accelerate the performance (Kaneko et al., 2012b).

# 244 3. Inference and Learning

#### 245 3.1. Inference

In inference, the maximum a posteriori (MAP) labeling of the random field is estimated:

$$\hat{\boldsymbol{y}} = \arg\max_{\boldsymbol{y} \in \mathcal{L}^N} P(\boldsymbol{y}|\boldsymbol{x}).$$
(12)

The exact distribution  $P(\boldsymbol{y}|\boldsymbol{x})$  for all the sets of labels  $\mathcal{L}^{N}$  is computationally intractable; however, the calculation cost is reduced to linear in the number of detected persons via the highly efficient approximation method described in Krähenbühl and Koltun (2011), because we define the pairwise potentials in our model as Gaussian kernels. The approximation method uses cross bilateral filtering techniques within a mean field approximation framework. The mean field approximation finds a product of independent marginals  $Q(\boldsymbol{y}) = \prod_i Q_i(y_i)$  close to  $P(\boldsymbol{y})$  in terms of minimizing the KLdivergence  $\mathbf{D}(Q||P)$  (Koller and Friedman, 2009). By considering the fixedpoint equations that hold at the stationary points of KL-divergence, the following iterative update equation is derived for  $Q_i(y_i = l)$  given  $Q_j(y_j)$  for all  $j \neq i$ :

$$Q_{i}(y_{i} = l) = \frac{1}{Z_{i}} \left\{ -\psi_{u}(y_{i}) - \sum_{l' \in \mathcal{L}} \mu(l, l') \sum_{j \neq i} k(\mathbf{f}_{i}, \mathbf{f}_{j}) Q_{j}(l') \right\}, \quad (13)$$

where  $Z_i$  normalizes the marginal at node *i*. A naive implementation of this 260 approximation has quadratic complexity in the number of variables N, when 261 calculating a message passing step:  $\tilde{Q}_i(l) = \sum_{j \neq i} k(f_i, f_j) Q_j(l)$ . However, in 262 Krähenbühl and Koltun (2011), it is shown that the message passing step can 263 be expressed as a convolution with a Gaussian kernel  $\tilde{Q}_i(l) = [G \otimes Q(l)](f_i) -$ 264  $Q_i(l)$ , where G is a Gaussian kernel and  $\otimes$  is the convolution operator, and it 265 is possible to leverage high-dimensional filtering, such as the permutohedral 266 lattice method (Adams et al., 2010). This reduces the calculation complexity 267 from quadratic to linear in the number of variables N. It should be noted 268 that Krähenbühl and Koltun (2011) report that the inference time is less 269 than one sec for several tens of thousands of variables. 270

#### 271 3.2. Learning

In learning, we use piecewise training (Sutton and McCallum, 2005). In piecewise training, the model is divided into several components, each of which is trained independently. Theoretically speaking, piecewise training minimizes an upper bound on the log partition function of the model. The

experimental results presented in Shotton et al. (2006) and Sutton and Mc-276 Callum (2005) show that piecewise training often performs comparably to 277 global training, even when joint full inference is used. In our model, the 278 unary potentials are trained first, using the contextual feature descriptors 279 described in Section 2.3. Next, we learn the kernel parameters  $w, \theta_1, \theta_2, \theta_3$ , 280 and  $\theta_4$  in the pairwise potentials. w can be found efficiently by exploiting 281 expectation maximization and high-dimensional filtering. However, it is not 282 easy to optimize the kernel widths  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$ , and  $\theta_4$  globally, due to their non 283 convexity on log-loss criteria. Therefore, we use a grid search on the training 284 set with cross-validation for all the kernel parameters. 285

# 286 4. Experiments

# 287 4.1. Datasets and Experimental Setup

We evaluated our model on the Collective Activity Dataset introduced in 288 Choi et al. (2009) and the dataset introduced in Choi and Savarese (2012). 289 Hereafter, we call the former the "Dataset I," and the latter the "Dataset 290 II." These datasets were considered appropriate for our evaluation, since 291 they include activities performed by multiple persons in a natural setting. 292 In most previous studies in the human action recognition field, the proposed 293 algorithms were evaluated on standard benchmark datasets such as the KTH 294 (Blank et al., 2005) and Weizmann (Schuldt et al., 2004) datasets. However, 295 these datasets include a single person performing a specific action, and the 296 video clips in the datasets were recorded in a controlled setting with a small 297 amount of camera motion and a clean background. The Hollywood human 298 action dataset (Laptev et al., 2008) and UT-Interaction dataset (Ryoo and 299

Aggarwal, 2010) are more challenging and contain actions performed by more than one actor; however, the actions (e.g., hand shaking, hugging) are not collective, but rather two persons perform one action together.

The Dataset I is composed of 44 short video clips, including five activity 303 classes: crossing, waiting, queuing, walking, and talking. The video clips in 304 the dataset were recorded using low resolution hand-held cameras under real-305 istic conditions, including camera shaking, background clutter, and transient 306 mutual occlusions of persons. Some video clips include multiple collective 307 activities in a scene or activity transition. All the persons in every tenth 308 frame of the videos are labeled with the ground truth: pose, activity, and 309 bounding box information. 310

The Dataset II is composed of 32 short video clips, including six activity 311 classes: gathering, talking, dismissal, walking together, chasing, and queuing. 312 The video clips in the dataset were recorded using a fixed camera in the out-313 doors. We omitted the activity queuing and used the remaining 30 short video 314 clips for evaluation, because the number of video clips containing queuing is 315 only two, and therefore, too small to evaluate using a leave-one-video-out 316 cross-validation scheme. It should be noted that each of the other activities 317 is contained in more than five video clips, respectively. While the pose label 318 and bounding box information are annotated per person, the activity label 319 is annotated per frame, since Choi and Savarese (2012) assumed there exists 320 only one collective activity in a certain time frame, and their goal was to clas-321 sify the activity per frame. However, this assumption is not always natural, 322 because some video clips contain multiple activities in a single frame, e.g., 323 one person is passing by, while the others are chasing. We therefore added 324

the walking alone activity class and re-annotated each person with six activities: gathering, talking, dismissal, walking together, chasing, and walking alone.

We evaluated our model in a way similar to that of Choi et al. (2009, 328 2011), Khamis et al. (2012a,b), Lan et al. (2010a), and Kaneko et al. (2012b). 329 For each dataset, we used the leave-one-video-out cross-validation scheme. 330 This means that when we classified activities in one video, all the other 331 videos in the dataset were used for training and validation. We report ac-332 tivity classification results on a per-person basis. It should be noted that, in 333 some previous studies (Lan et al., 2010b; Choi and Savarese, 2012), it was 334 assumed that there exists only one collective activity in a certain time frame, 335 and activity classification results were reported on a per-frame basis. This 336 experimental protocol is inadequate to evaluate our model, because our goal 337 was to segment collective activities in a scene where multiple groups exist, 338 as shown in Fig. 2. 339

The experiments were conducted on an Intel Core i7 processor clocked at 2.2 GHz. The calculation complexity of inference using our fully connected model was O(N), where N (the number of detected persons) was typically in the order of tens, hundreds, or thousands, and inference time was less than 0.1 sec per video sequence for both the datasets. It should be noted that, in Choi and Savarese (2012), the classification and target association take about 1 min per video sequence, given tracklets and observations.

#### 347 4.2. Evaluation of Graph Structures

In order to evaluate the performance of the proposed model comprehensively, we compared it with several baseline models. The first baselines use



Figure 3: A structure of person-person interactions in each graph model. Each node represents a person in a video. Dashed lines represent that the interaction strength between persons is altered according to the degree of their similarity, while solid lines represent that the interaction strength between persons is constant, regardless of the degree of their similarity. (a) The unary only model (no connection between any pair of nodes); (b) nodes are connected in a frame; (c) nodes are connected in adjacent frames; (d) and (e) all the nodes are connected in a video.

various ways of setting the structures of the person-person interactions. The 350 structures that we considered are shown in Fig. 3(a)-(c), including (a) the 351 Unary Only model (no pairwise connection); (b) the graph obtained by con-352 necting persons in a frame (the *Connected Per Frame* model); (c) the graph 353 obtained by connecting persons in adjacent frames (the Adjacently Connected 354 model). It should be noted that in our proposed Fully Connected model, all 355 the persons in a video are connected, as shown in Fig. 3(d). In the three 356 models shown in Fig. 3(b)-(d), pairwise potentials are defined as Gaussian 357 edge potentials, and therefore, the interaction strength between persons is 358 altered according to the degree of their similarity. 359

In order to evaluate the Gaussian edge potentials, we also compared the proposed model with the other baseline (which we call *Simple Fully Connected*), fully connected CRFs with constant edge potentials, as shown in Fig. 3(e). In the model, the pairwise potential in equation (2) is defined as  $\phi_p(y_i, y_j) = w\mu(y_i, y_j)$ . Unlike in the proposed model, the interaction strength between persons in this model is constant over a video. In order to
analyze the efficacy of the graph structure models, we used the AC descriptor (Lan et al., 2010a) as the baseline, and calculated the unary potential in
equation (3) based on this descriptor.

Table 1: Comparison of activity classification accuracy levels for different graph structures. Each graph structure has a different structure of person-person interactions or different interaction strength between persons. The details are illustrated in Fig. 3.

Method/Dataset	Dataset I	Dataset II	
	(Choi et al., 2009)	(Choi and Savarese, 2012)	
Unary Only	67.4%	56.2%	
Connected Per Frame	68.6%	56.7%	
Adjacency Connected	69.6%	57.6%	
Fully Connected	72.2%	$\mathbf{59.3\%}$	
Simple Fully Connected	68.7%	55.5%	

Quantitative Results. The comparison results of our approach and the base-369 lines are summarized in Table 1. The results shows that our model (Fully 370 *Connected*) outperforms the baselines on the two datasets. The *Simple Fully* 371 *Connected* model has a fully connected graph structure, as does the *Fully* 372 *Connected* model; however, the former is significantly outperformed by the 373 latter, since the former exploits the long-range interaction but cannot deal 374 with multiple collective activities and activity transition. The confusion ma-375 trices of our model and the baseline (the Unary Only model) are illustrated 376 in Fig 4. These confusion matrices show that our proposed Fully Connected 377 model achieves a significant improvement over the baseline. It should be 378



Figure 4: Confusion matrices for activity classification on the two datasets. (a) and (b) show the confusion matrices for collective activity using the baseline (the *Unary Only* model) and proposed method on the Dataset I (Choi et al., 2009), respectively. (c) and (d) compare the two methods on the Dataset II (Choi and Savarese, 2012). In both cases, our fully connected model significantly outperforms the baseline method. In a confusion matrix, rows represent ground truth and columns represent prediction. Each row is normalized to sum to 1.

noted that walking vs crossing for the Dataset I is still ambiguous in our 379 model, because whether a person is *walking* or *crossing* often depends not on 380 the person-person interactions but on the environmental setting, such as a 381 sidewalk or a pedestrian crossing. In both our model and the baseline, walk-382 ing alone on the Dataset II is confused with dismissal and walking together, 383 because *walking alone* is a single person activity and the neighbors of the 384 person have a variety of appearances, e.g., one person passing by other peo-385 ple often seems like dismissal. However, the number of confusions is reduced 386 using our model. 387

Qualitative Results. We also visualize the qualitative results using different structures of person-person interactions in Figs. 5–7. Fig. 5 shows typical success and failure examples, Fig. 6 shows examples of the scenes where multiple groups exist, and Fig. 7 shows examples of the case where the activities <sup>392</sup> transition from one to another.

Fig. 5 shows typical success and failure examples for the Dataset I. Fig. 5(a)-393 (b) demonstrates that the classification result achieved by the Unary Only 394 model is often noisy and inconsistent, since the activity of each person is 395 classified independently in the model. In the *Connected Per Frame* model 396 and the Adjacently Connected model, collective activities are optimized in 397 a frame and consecutive frames, respectively; however, misclassification of 398 some persons in the duration of the video clip often causes misclassification 399 of others. In contrast, in our *Fully Connected* model, the misclassification 400 is fixed, and we obtain the temporal and spatial consistency of the activity 401 in a group by leveraging the cues over the multi-scale. However, the wrong 402 classification in the Unary Only model causes any graph structure to classify 403 everyone incorrectly, as shown in Fig. 5(c). 404

Fig. 6 visualizes examples of the scenes where multiple groups exist in 405 the Dataset I. Fig. 6(a)-(b) shows that the Fully Connected model yields 406 misclassification errors, but achieves consistency in each group, while the 407 other models fail to obtain consistency. This is because the Connected Per 408 Frame model and the Adjacently Connected model are likely to be influenced 409 by the misclassification of some persons in the same frame and consecutive 410 frames, respectively. However, the high level of confidence in the wrong label 411 in the Unary Only model causes any structure to fail to obtain consistency, 412 as shown in Fig. 6(c). 413

Fig. 7 visualizes examples of the case where the activities transition from one to another in the Dataset II. In particular, the activity transitions from *gathering* to *talking*. Since the point at which *gathering* changes to *talking* is ambiguous, that of the *Fully Connected* model is not completely consistent
with that of the ground truth; however, only our model succeeds in recognizing the transition between the two activities, in contrast to any other
model.



Figure 5: (Best viewed in color) Visualization of typical success and failure examples for the Dataset I (Choi et al., 2009) using different structures of person-person interactions. The labels C (magenta), S (blue), Q (cyan), W (red), T (green), and NA (white) indicate crossing, waiting, queuing, walking, talking, and not assigned, respectively. The first row shows the ground truth; the second row shows results using the Unary Only model; the third row shows results using the Connected Per Frame model; the fourth row shows results using the Adjacently Connected model; the fifth row shows results using our proposed Fully Connected model. The first three columns (a) represent results in the consecutive frames in the same video clip, and the last two columns (b)-(c) represent results in other video clips. In the first four columns (a)-(b), temporary misclassification in the Unary Only model is fixed in the Fully Connected model. In the fifth column (c), any structure fails to classify collective activities, affected by the wide and continuous misclassification in the Unary Only model.



Figure 6: (Best viewed in color) Visualization of examples of the scenes where multiple groups exist in the Dataset I (Choi et al., 2009). The labels C (magenta), S (blue), Q (cyan), W (red), T (green), and NA (white) indicate crossing, waiting, queuing, walking, talking, and not assigned, respectively. The first row shows the ground truth; the second row shows results using the Unary Only model; the third row shows results using the Connected Per Frame model; the fourth row shows results using the Adjacently Connected model. The first three columns (a) represent results in the consecutive frames in the same video clip, and the last two columns (b)-(c) represent results in other video clips. In the first four columns (a)-(b), the Fully Connected model achieves consistency in a group, although the other models fail to obtain consistency in some scenes. In the fifth column (c), the high level of confidence in the wrong label in the Unary Only model causes any structure to fail to obtain consistency.



Figure 7: (Best viewed in color) Visualization of examples of the case where the activities transition from one to another in the Dataset II (Choi and Savarese, 2012). In particular, the activity transitions from *gathering* to *talking*. The labels G (magenta), T (green), D (cyan), and W (red) indicate *gathering*, *talking*, *dismissal*, and *walking together*, respectively. The first row shows the ground truth; the second row shows results using the *Unary Only* model; the third row shows results using the *Connected Per Frame* model; the fourth row shows results using the *Adjacently Connected* model; the fifth row shows results using our proposed *Fully Connected* model. Since the point at which *gathering* changes to *talking* is ambiguous, that of the *Fully Connected* model is not completely consistent with that of the ground truth; however, only our model succeeds in recognizing the transition between the two activities, in contrast to any other model.

#### 421 4.3. Comparison with State-of-the-Art Methods

We also compared our method with recent methods (Choi et al., 2009, 422 2011; Khamis et al., 2012a,b; Lan et al., 2010a; Kaneko et al., 2012b). So 423 that the comparison would be fair, we used the same leave-one-video-out 424 scheme described in the studies; we report activity classification results on 425 a per-person basis. In order to evaluate our fully connected model in com-426 parison with state-of-the art models, we calculated the unary potentials in 427 equation (3) using not only the AC descriptor (Lan et al., 2010a) but also 428 the combination of the AC and RAC descriptors (AC-RAC) (Kaneko et al., 429 2012b). 430

*Results.* The comparison results of the activity classification accuracy levels
for different methods for the Dataset I and Dataset II are summarized in
Table 2. We analyzed our results mainly on the Dataset I, because in no
previous study a model was evaluated using the same protocol as ours in the
Dataset II.

In the models using the AC descriptor as the baseline, our model (AC +436 FC-CRF) outperforms the other models: AC + Frame Cues (Khamis et al., 437 2012a), AC + Track Cues (Khamis et al., 2012b) and AC + Frame/Track438 Cues (Khamis et al., 2012a), although our baseline (AC) is inferior to those 439 of (Khamis et al., 2012a,b). It should be noted that AC + Frame/Track440 Cues (Khamis et al., 2012a) connects only people with overlapping bounding 441 boxes in consecutive frames to improve the inference speed, while our model 442 can solve the fully connected model in less than 0.1 sec. AC + FC-CRF also 443 outperforms STV + MC (Choi et al., 2009) and RSTV + MRF (Choi et al., 444 2011). It should be noted that these models (Choi et al., 2009, 2011) em-44

<sup>446</sup> ploy additional trajectory information of each person to obtain consistency.
<sup>447</sup> However, it is not easy to obtain the correct trajectory in a scene where
<sup>448</sup> transient mutual occlusions of persons exist. Mistakes produced during the
<sup>449</sup> tracking step can influence the performance of the model during recognition.
<sup>450</sup> Khamis et al. (2012a,b) also did not use the trajectory, and instead solved
<sup>451</sup> the identity maintenance problem and action recognition simultaneously.

We also show the results obtained using our previously proposed descrip-452 tor (AC-RAC) (Kaneko et al., 2012b) and integration of our model and the 453 descriptor (AC-RAC + FC-CRF). In the two datasets, AC-RAC outperforms 454 the other feature descriptors, and the integration of our model and the de-455 scriptor (AC-RAC + FC-CRF) achieves the best performance. Choi and 456 Savarese (2012) reported an accuracy (74.4%) that is competitive with that 457 achieved by our method (74.7%) on the Dataset I. However, since their ac-458 curacy level was achieved using a different experimental protocol, i.e., they 459 assigned the per-scene collective activity labels that they obtained with four-460 fold experiments to each individual, it is not directly comparable to the 461 accuracy levels listed in Table 2. 462

In the Dataset II, no previous method was directly comparable to our 463 method, because, in no previous study, a model was evaluated on a *per-person* 464 basis like ours, while Choi and Savarese (2012) evaluated their model on a 465 per-frame basis, and reported an accuracy (74.3% mean-per-class using base-466 line method and 79.2% mean-per-class using their proposed method). The 467 direct comparison between our method and their method using the same ex-468 perimental protocol would be an interesting topic, however, it would be hard 469 to perform the direct comparison, because the assumptions of our and their 470

models are basically different and it is required to reformulate the problems 471 significantly. Our model aims to recognize an activity in each group even 472 in the scene where multiple activity groups exist, while their model aims to 473 recognize one main collective activity per frame under the assumption that 474 there is only one collective activity per scene per time stamp. It should be 475 noted that our and their model require the different annotated information 476 in training process. Our model requires the activity label annotated per per-477 son, while their model requires the activity label annotated per frame. In the 478 scene where some activity groups exist, they assigned each frame into one 479 activity label by taking the majority of activities of persons in that frame. 480

Table 2: Comparison of activity classification accuracy levels for different methods on the Dataset I (Choi et al., 2009) and Dataset II (Choi and Savarese, 2012). The methods without a plus sign (+) do not use graph structure, s.t., unary only models, such as shown in Fig. 3(a), while the methods with a plus sign (+) use a graph structure, such as shown in Fig. 3(b)–(e). The last four rows show the results obtained using our implementation. In the two datasets, our previously proposed descriptor (*AC-RAC*) (Kaneko et al., 2012b) outperforms the other feature descriptors, and the integration of our model and the descriptor (*AC-RAC* + *FC-CRF*) achieves the best performance.

Method/Dataset	Dataset I	Dataset II
	(Choi et al., 2009)	(Choi and Savarese, 2012)
AC (Lan et al., 2010a)	68.2%	-
STV (Choi et al., $2009$ )	64.3%	-
STV + MC (Choi et al., 2009)	65.9%	-
RSTV (Choi et al., 2011)	67.2%	-
RSTV + MRF (Choi et al., 2011)	70.9%	-
AC (Khamis et al., 2012b)	68.8%	-
AC + Track Cues (Khamis et al., 2012b)	70.9%	-
AC + Frame Cues (Khamis et al., 2012a)	70.7%	-
AC + Frame/Track Cues (Khamis et al., 2012a)	72.0%	-
AC	67.4%	56.2%
AC + FC-CRF	72.2%	59.3%
AC-RAC	71.9%	66.4%
AC-RAC + FC-CRF	74.7%	70.7%

#### 481 5. Conclusion and Discussion

In this paper, we described a method for consistent collective activity recognition that uses fully connected CRFs, which consider all the interactions among the people in a video clip and alter the interaction strength according to the degree of their similarity. Our model leverages various features, such as position, size, motion, and time sequence, over a "multi-scale"

in a single unified model, in order to allow various types, sizes, and shapes 487 of groups to be treated. The intractability of the fully connected model is 488 overcome by describing the pairwise potentials as Gaussian kernels and using 489 a highly efficient approximation method. Our experimental results showed 490 that our model is robust against temporary misclassification and able to deal 491 with multiple activities in a scene and activity transition. Quantitative eval-492 uation on the two publicly available datasets showed that our fully connected 493 model outperforms state-of-the art models, as well as other graph structures, 494 such as the unary only model and the adjacently connected model. 495

Future directions include investigation of other useful contexts for collective activity recognition, such as environmental settings; extension to online processing for online applications; and research on efficient learning techniques as a substitute for grid search to optimize the kernel parameters efficiently.

# 501 References

- Adams, A., Baek, J., Davis, M. A., 2010. Fast high-dimensional filtering using the permutohedral lattice. In: Computer Graphics Forum.
- Amer, M. R., Todorovic, S., 2011. A chains model for localizing participants
  of group activities in videos. In: International Conference on Computer
  Vision (ICCV).
- <sup>507</sup> Blank, M., Gorelick, L., Shechtman, E., Irani, M., Basri, R., 2005. Actions
  <sup>508</sup> as space-time shapes. In: International Conference on Computer Vision
  <sup>509</sup> (ICCV).

- <sup>510</sup> Boykov, Y. Y., Jolly, M.-P., 2001. Interactive graph cuts for optimal bound<sup>511</sup> ary and region segmentation of objects in n-d images. In: International
  <sup>512</sup> Conference on Computer Vision (ICCV).
- <sup>513</sup> Choi, W., Savarese, S., 2012. A unified framework for multi-target tracking
  <sup>514</sup> and collective activity recognition. In: European Conference on Computer
  <sup>515</sup> Vision (ECCV).
- <sup>516</sup> Choi, W., Shahid, K., Savarese, S., 2009. What are they doing?: Collective
  <sup>517</sup> activity classification using spatio-temporal relationship among people. In:
  <sup>518</sup> International Workshop on Visual Surveillance (VS).
- <sup>519</sup> Choi, W., Shahid, K., Savarese, S., 2011. Learning context for collective
  <sup>520</sup> activity recognition. In: IEEE Conference on Computer Vision and Pattern
  <sup>521</sup> Recognition (CVPR).
- Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human detection. In: IEEE Conference on Computer Vision and Pattern Recognition
  (CVPR).
- Fan, R. E., Chang, K. W., Hsieh, C. J., Wang, X. R., Lin, C. J., 2008.
  LIBLINEAR: A library for large linear classification. Journal of Machine
  Learning Research (JMLR) 9, 1871–1874.
- Felzenszwalb, P., McAllester, D., Ramanan, D., 2008. A discriminatively
  trained, multiscale, deformable part model. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- <sup>531</sup> Hatef, M., Duin, R. P., Matas, J., 1998. On combining classifiers. IEEE

- Transaction on Pattern Recognition and Machine Intelligence (PAMI) 20,
  226–239.
- He, X., Zemel, R. S., Carreira-Perpiñán, M. A., 2004. Multiscale conditional
  random fields for image labeling. In: IEEE Conference on Computer Vision
  and Pattern Recognition (CVPR).
- Kaneko, T., Shimosaka, M., Odashima, S., Fukui, R., Sato, T., 2012a. Consistent collective activity recognition with fully connected CRFs. In: International Conference on Pattern Recognition (ICPR).
- Kaneko, T., Shimosaka, M., Odashima, S., Fukui, R., Sato, T., 2012b. Viewpoint invariant collective activity recognition with relative action context.
  In: Workshop on Action Recogniton and Pose Estimation in Still Images
  (APSI).
- Khamis, S., Morariu, V. I., Davis, L. S., 2012a. Combining per-frame and pertrack cues for multi-person action recognition. In: European Conference
  on Computer Vision (ECCV).
- Khamis, S., Morariu, V. I., Davis, L. S., 2012b. A flow model for joint action
  recognition and identity maintenance. In: IEEE Conference on Computer
  Vision and Pattern Recognition (CVPR).
- Koller, D., Friedman, N., 2009. Probabilistic Graphical Models: Principles
   and Techniques. MIT Press.
- 552 Krähenbühl, P., Koltun, V., 2011. Efficient inference in fully connected CRFs
- with Gaussian edge potentials. In: Advances in Neural Information Processing Systems (NIPS).

- Lafferty, J., McCallum, A., Pereira, F., 2001. Conditional random fields:
  Probabilistic models for segmenting and labeling sequence data. In: International Conference on Machine Learning (ICML).
- Lan, T., Wang, Y., Mori, G., 2010a. Retrieving actions in group contexts.
  In: International Workshop on Sign Gesture Activity (SGA).
- Lan, T., Wang, Y., Yang, W., Mori, G., 2010b. Beyond actions: Discriminative models for contextual group activities. In: Advances in Neural
  Information Processing Systems (NIPS).
- Laptev, I., Marszalek, M., Schmid, C., Rozenfeld, B., 2008. Learning realistic
  human actions from movies. In: IEEE Conference on Computer Vision and
  Pattern Recognition (CVPR).
- Niebles, J. C., Wang, H., Fei-Fei, L., 2006. Unsupervised learning of human
  action categories using spatial-temporal words. In: British Machine Vision
  Conference (BMVC).
- <sup>569</sup> Ryoo, M. S., Aggarwal, J. K., 2010. UT-Interaction Dataset, ICPR
  <sup>570</sup> contest on Semantic Description of Human Activities (SDHA).
  <sup>571</sup> http://cvrc.ece.utexas.edu/SDHA2010/Human\_Interaction.html.
- Schuldt, C., Laptev, I., Caputo, B., 2004. Recognizing human actions: A
  local SVM approach. In: International Conference on Pattern Recognition
  (ICPR).
- Shotton, J., Winn, J., Rother, C., Criminisi, A., 2006. Textonboost: Joint
  appearance, shape and context modeling for multi-class object recognition
  and segmentation. In: European Conference on Computer Vision (ECCV).

- Sun, D., Roth, S., Black, M. J., 2010. Secrets of optical flow estimation and
  their principles. In: IEEE Conference on Computer Vision and Pattern
  Recognition (CVPR).
- <sup>581</sup> Sutton, C., McCallum, A., 2005. Piecewise training for undirected models.
- In: Conference on Uncertainty in Artificial Intelligence (UAI).