# Deep Structured Models For Group Activity Recognition

Zhiwei Deng<sup>1</sup> zhiweid@sfu.ca Mengyao Zhai<sup>1</sup> mzhai@sfu.ca Lei Chen<sup>1</sup> chenleic@sfu.ca Yuhao Liu<sup>1</sup> yla305@sfu.ca Srikanth Muralidharan<sup>1</sup> smuralid@sfu.ca Mehrsan Javan Roshtkhari<sup>2</sup> mehrsan@sportlogiq.com Greg Mori<sup>1</sup> mori@cs.sfu.ca

- <sup>1</sup> School of Computing Science Simon Fraser University Burnaby, BC, Canada
- <sup>2</sup> SportLogiq Inc. Montreal, QC, Canada

#### Abstract

This paper presents a deep neural-network-based hierarchical graphical model for individual and group activity recognition in surveillance scenes. As the first step, deep networks are used to recognize activities of individual people in a scene. Then, a neuralnetwork-based hierarchical graphical model refines the predicted labels for each activity by considering dependencies between different classes. Similar to the inference mechanism in a probabilistic graphical model, the refinement step mimics a message-passing encoded into a deep neural network architecture. We show that this approach can be effective in group activity recognition and the deep graphical model improving recognition rates over baseline methods.

# **1** Introduction

Event understanding in videos is a key element of computer vision systems in the context of visual surveillance, human-computer interaction, sports interpretation, and video search and retrieval. Therefore events, activities, and interactions must be represented in such a way that retains all of the important visual information in a compact and rich structure. Accurate detection and recognition of atomic actions of each individual person in a video is the primary component of such a system, and also the most important, as it affects the performance of the whole system significantly. Although there are many methods to determine human actions in uncontrolled environments, this task remains a challenging computer vision problem, and robust solutions would open up many useful applications. The standard

© 2015. The copyright of this document resides with its authors. It may be distributed unchanged freely in print or electronic forms.



Figure 1: Recognizing individual and group activities in a deep network. Individual action labels are predicted via CNNs. Next, these are refined through a message passing neural network which considers the dependencies between the predicted labels.

and yet state-of-the-art pipeline for activity recognition and interaction description consists of extracting hand-crafted local feature descriptors either densely or at a sparse set of interest points (e.g., HOG, MBH, ...) in the context of a Bag of Words model [22]. These are then used as the input either to a discriminative or a generative model. In recent years, it has been shown that deep learning techniques can achieve state-of-the-art results for a variety of computer vision tasks including action recognition [12].

On the other hand, understanding of complex visual events in a scene requires exploitation of richer information rather than individual atomic activities, such as recognizing local pairwise and global relationships in a social context and interaction between individuals and/or objects [5, 13, 13, 13, 13, 12]. This complex scene description remains an open and challenging task. It shares all of the difficulties of action recognition, interaction modeling<sup>1</sup>, and social event description. Formulating this problem within the probabilistic graphical models framework provides a natural and powerful means to incorporate the hierarchical structure of group activities and interactions [12], [13]. Given the fact that deep neural networks can achieve very competitive results on the single person activity recognition tasks, they can, produce better results when they are combined with other methods, *e.g.* graphical models, in order to capture the dependencies between the variables of interest [21]. Following a similar idea of incorporating spatial dependency between variables into the deep neural network in a joint-training process presented [21], here we focus on learning interactions and group activities in a surveillance scene by employing a graphical model in a deep neural network paradigm.

In this paper, our main goal is to address the problem of *group activity understanding* and *scene classification* in complex surveillance videos using a deep learning framework. More specifically, we are focused on learning individual activities and describing the scene

<sup>&</sup>lt;sup>1</sup>The term "interaction" refers to any kind of interaction between humans, and humans and objects that are present in the scene, rather than activities which are performed by a single subject.



Figure 2: A schematic overview of our message passing CNN framework. Given an image and the detected bounding boxes around each person, our model predicts scores for individual actions and the group activity. The predicted labels are then refined by applying a belief propagation-like neural network. This network considers the dependencies between individual actions, body poses, and the group activity. The model learns the message passing parameters and performs inference and learning in unified framework using back-propagation.

simultaneously while considering the pair-wise interactions between individuals and their global relationship in the scene. This is achieved by combining a Convolutional Neural Network (CNN) with a probabilistic graphical model as additional layers in a deep neural network architecture into a unified learning framework. The probabilistic graphical models can be seen as a refining process for predicting class labels by considering dependencies between individual actions, body poses, and group activities. The probabilistic graphical model is modeled by a multi-step message passing neural network and the predicted label refinement is carried out through belief propagation layers in the neural network. Figure 1 depicts an overview of our approach for label refinement. Experimental results show the effectiveness of our algorithm in both activity recognition and scene classification.

# 2 Previous Work

The analysis of human activities is an active area of research. Decades of research on this topic have produced a diverse set of approaches and a rich collection of activity recognition algorithms. Readers can refer to recent surveys such as Poppe [16] and Weinland et al. [23] for a review. Many approaches concentrate on an activity performed by a single person, including state of the art deep learning approaches [11], [19].

In the context of scene classification and group activity understanding, many approaches use a hierarchical representation of activities and interactions for collective activity recognition [ $\square$ ]. They have been focused on capturing spatio-temporal relationships between visual cues either by imposing a richer feature descriptor which accounts for context [ $\square$ ,  $\square$ ] or a context-aware inference mechanism [ $\square$ ,  $\square$ ]. Hierarchical graphical models [ $\square$ ,  $\square$ ,  $\square$ ,  $\square$ ],

AND-OR graphs [**D**, **D**], and dynamic Bayesian networks [**D**] are among the representative approaches for group activity recognition.

In traditional approaches, local hand-crafted features/descriptors have been employed to recognize atomic activities. Recently, it has been shown that the use of deep neural networks can by itself outperform other algorithms for atomic activity recognition. However, no prior art in the CNN-based video description used activities and scene information jointly in a unified graphical representation for scene classification. Therefore, the main objective of this research is to develop a system for activity recognition and scene classification which simultaneously uses the action and scene labels in a neural network-based graphical model to refine the predicted labels via a multiple-step message passing procedure.

More closely related to our approach are the ones combining graphical models with convolutional neural networks [B, [20]]. In [20], a one step message passing is implemented as a convolution operation in order to incorporate spatial relationship between local detection responses for human body pose estimation. In another study, Deng *et al.* [B] propose an interesting solution to improve label prediction in large scale classification by considering relations between the predicted class labels. They employ a probabilistic graphical model with hard constraints on the labels on top of a neural network in a joint training process. In essence, our proposed algorithm follows a similar idea of considering dependencies between predicted labels for the actions, group activities, and the scene label to solve the group activity recognition problem. Here we focus on incorporating those dependencies by implementing the label refinement process via an inter-activity neural network, as shown in Figure 2. The network learns the message passing procedure and performs inference and learning in unified framework using the back-propagation algorithm.

### 3 Model

Considering the architecture of our proposed structured label refinement algorithm for group activity understanding (see Figure 2), the key part of the algorithm is a multi-step message passing neural network. In this section, we describe how to combine neural networks and graphical models by mimicking a message passing algorithm and how to carry out the training procedure.

#### 3.1 Graphical Models in a Neural Network

Graphical models provide a natural way to hierarchically model group activities and capture the semantic dependencies between the group and individual activities [12]. A graphical model defines a joint distribution over states of a set of nodes. For instance, one can use a factor graph, in which each  $\phi_i$  corresponds to a factor over a set of related variable nodes  $x_i$  and  $y_i$ , and models interactions between those nodes in a log-linear fashion:

$$P(X,Y) \propto \prod_{i} \phi_i(x_i, y_i) \propto exp(\sum_k w_k f_k(x, y))$$
(1)

where X represents the inputs and Y is the predicted labels, with a weighted  $(w_k)$  feature functions  $f_k$ .

In order to do the inference in a graphical model, belief propagation is often adopted as a way to infer states or probabilities of the variables. In the belief propagation algorithm, each step of message passing involves two parts. At first the relevant information from the



Figure 3: Weight sharing scheme in a neural network. We use a sparsely connected layer to represent message passing between variable and factor nodes. Each factor node only connects to the relevant nodes. The factor nodes of same type share the same template of the parameters. For example, the first two factor nodes (the left and the middle one) have the same type and hence, share the same set of parameters which are the information from the scene1, action1 and pose1. The third factor node (the right one) adopts another set of weights.

connected nodes to a factor node are collected. Those messages are the passed to the variable nodes by marginalizing over states of irrelevant variables.

Following this idea, we mimic the message passing process by representing each combination of states as a neuron in a neural network, denoted as a "factor neuron". While normal message passing calculates dependencies rigidly, a factor neuron can be used to learn and predict dependencies between states and pass messages to the variable nodes. In the setting of neural networks, this dependency representation becomes more flexible and can adopt varied types of neurons (linear, ReLU, Sigmoid, etc.). Moreover, by integrating graphical models into a neural network, the formulation of a graphical model allows for parameter sharing in the neural network. Parameter sharing not only reduces the number of free parameters to learn, but also accounts for the semantic similarities between factor neurons. Figure 3 shows the parameter sharing scheme for different factor neurons.

### 3.2 Message Passing CNN Architecture for Group Activity

Representing group activities and individual activities as a hierarchical graphical model has proven to be a successful strategy  $[\Box, \Box, \Box\Box]$ . We adopt a similar structured model that considers group activity, individual activities, and group-individual interactions together. We introduce a new message passing convolutional neural network framework as shown in Figure 2. The model has two main parts: (i) a set of fine-tuned CNNs that produce a *scene score* for an image, and *action scores* and *pose scores* for each individual person in that image; and (ii) a message passing neural network which captures the dependencies between activities, poses, and scene labels.

Given an image *I*, and a set of bounding boxed for detected persons  $\{I_1, I_2, ..., I_M\}^2$ , the first part of our model generates raw scores of scene. In addition, it produces the raw scores for the actions and poses of each of the *M* individuals in the image  $\{I_i\}_{i=1}^M$ . This is done by applying fine-tuned CNNs on the image and the detected bounding boxes. A soft-max normalization is then applied for each scene, activity, and pose score in order to produce the raw scores.

The second part of our algorithm which does the label refinement takes the those raw scores as the imput. In our graphical model, outputs from CNNs correspond to unary poten-

 $<sup>^{2}</sup>$ It is assumed that the bounding box around each person is known. Those bounding boxes are obtained by applying a person detector on each image as a pre-processing step.

tials. The scene-level, and per-person action and pose-level unary potentials for the image I are represented by  $\mathbf{s}^{(k)}(I)$ ,  $\mathbf{a}^k(I_m)$ , and  $\mathbf{r}^{(k)}(I_m)$  respectively. The superscript  ${}^{(k)}$  is the index of message passing step. We use G to denote all group activity labels, H to represent all the action labels and Z to denote all the pose labels. Then the group activity in one scene can be represented as  $g_I$ ,  $\{h_{I_1}, h_{I_2}, ..., h_{I_M}\}$ ,  $\{z_{I_1}, z_{I_2}, ..., z_{I_M}\}$  where  $g_I \in G$  is the group activity label for image I,  $h_{I_i}$  and  $z_{I_i}$  are action labels and pose labels for a person  $I_m$ .

Note that for training, the scene, action, and pose CNN models in first part of our algorithm are fine-tuned from an AlexNet architecture pre-trained using the ImageNet data. The architecture is similar to the one proposed by  $[\square]$  for object classification with some minor differences, *e.g.* pooling which is done before the normalization. The network consists of five convolutional layers followed by two fully connected layers, and a softmax layer that outputs individual class scores. We use the softmax loss, stochastic gradient descent and dropout regularization to train these three CNNs.

In the second part of our algorithm, we use the method described in Sec. 3.1 to mimic the message passing in a hierarchical graphical model for group activity recognition in a scene. This stage can contain several steps of message passing. In each step, there are two types of passes: from outputs of step k - 1 to factor layer and from factor layer to k step outputs. In the kth message passing step, the first pass computes dependencies between states. The inputs to the kth step message passing are

$$\{s_1^{(k-1)}(I), \dots, s_{|G|}^{(k-1)}(I), a_1^{(k-1)}(I_1), \dots, a_{|H|}^{(k-1)}(I_M), r_1^{(k-1)}(I_1), \dots, r_{|Z|}^{(k-1)}(I_M)\}$$
(2)

where  $s_g^{(k-1)}(I)$  is the scene score of image *I* for label *g*,  $a_h^{(k-1)}(I_m)$  is the action score of person  $I_m$  for label *h* and  $r_z^{(k-1)}(I_m)$  is the pose score of person  $I_m$  for label *z*. In the factor layer, the action, pose and scene interaction are calculated as:

$$\phi_j(s_g^{(k-1)}(I), a_h^{(k-1)}(I_m), r_z^{(k-1)}(I_m))) = \alpha_{g,h,z}[s_g^{(k-1)}(I), a_h^{(k-1)}(I_m), r_z^{(k-1)}(I_m))]^T$$
(3)

where  $\alpha_{g,h,z}$  is a 3-d parameter template for combination of scene *g*, action *h* and pose *z*. Similarly, pose interactions for all people in the scene are calculated as:

$$\psi_j(s_g^{(k-1)}(I), \mathbf{r}) = \beta_{tg}[s_g^{(k-1)}(I), \mathbf{r}]^T$$
(4)

where **r** is all output nodes for all people, *t* is the factor neuron index for scene *g*. *T* latent factor neurons are used for a scene *g*. Note that parameters  $\alpha$  and  $\beta$  are shared within factors that have the same semantic meaning. For the output of *k*th step message passing, the score for the scene label to be *g* can be defined as:

$$s_{g}^{(k)}(I) = s_{g}^{(k-1)}(I) + \sum_{j \in \varepsilon_{1}^{s}} w_{ij}\phi_{j}(s_{g}^{(k-1)}(I), \mathbf{a}, \mathbf{r}; \alpha)) + \sum_{j \in \varepsilon_{2}^{s}} w_{ij}\psi_{j}(s_{g}^{(k-1)}(I), \mathbf{r}; \beta)$$
(5)

where  $\varepsilon_1^s$  and  $\varepsilon_2^s$  are the set of factor nodes that connected with scene *g* in first factor component(scene-action-pose factor) and second factor component (pose-global factor) respectively. Similarly, we also define action and pose scores after the *k*th message passing step as:

$$a_{h}^{(k)}(I_{m}) = a_{h}^{(k-1)}(I_{m}) + \sum_{j \in \varepsilon_{1}^{a}} w_{ij}\phi_{j}(a_{h}^{(k-1)}(I_{m}), \mathbf{s}, \mathbf{r}; \boldsymbol{\alpha})$$
(6)

$$r_{z}^{(k)}(I_{m}) = r_{z}^{(k-1)}(I_{m}) + \sum_{j \in \mathcal{E}_{1}^{r}} w_{ij}\phi_{j}(r_{z}^{(k-1)}(I_{m}), \mathbf{a}, \mathbf{s}; \boldsymbol{\alpha}) + \sum_{j \in \mathcal{E}_{2}^{r}} w_{ij}\psi_{j}(r_{z}^{(k-1)}(I_{m}), \mathbf{r}; \boldsymbol{\beta})$$
(7)

where  $\varepsilon = \{\varepsilon_1^s, \varepsilon_2^s, \varepsilon_1^a, \varepsilon_1^r, \varepsilon_2^r\}$  are connection configurations in the pass from factor neurons to output neurons. These connections are simply the reverse of the configurations in the first pass, from input to factors. The model parameters  $\{\mathbf{W}, \alpha, \beta\}$  are weights on the edges of the neural network. Parameter **W** represents the concatenation of weights connected from factor layers to output layer (second pass), while  $\alpha, \beta$  represent weights from the input layer of the  $k^{th}$  message passing to factor layers (first pass).

#### **3.2.1** Components in the Factor Layers

This section summarizes and explains all different components of our model, which are as follows:

**Unary component:** In the message passing model, the unary component corresponds to group activity scores for an image *I*, action and pose scores for each person  $I_m$  in frame *I*, represented as  $s_g^{(k-1)}(I)$ ,  $a_h^{(k-1)}(I_m)$  and  $r_z^{(k-1)}(I_m)$  respectively. These scores are acquired from the previous step of message passing and are directly added to the output of the next message passing step.

**Group activity-action-pose factor layer**  $\phi$ : A group's activity is strongly correlated to the participating individuals' actions. This component for the model is used to measure the compatibility between individuals and groups. An individual's activity can be described by both pose and action, and we use this ternary scene-pose-action factor layer to capture dependencies between a person's fine-grained action (e.g. talking facing front-left) and the scene label for a group of people. Note that in this factor layer we used the weight sharing scheme mentioned in Sec. 3.1 to mimic the belief propagation.

**Poses-all factor layer**  $\psi$ : Pose information is very important in understanding a group activity. For example, when all people are looking in the same direction, there is a high probability that it's a queueing scene. This component captures this global pose information for a scene. Instead of naively enumerate all combination of poses for all people, we exploit the sparsity of truly useful and frequent patterns, and simply use *T* factor nodes for one scene label. In our experiments, we simply set *T* to be 10.

#### 3.3 Multi Step Message Passing CNN Training

The steps of message passing depends on the structure of graphical model. In general, graphical models with loops or large number of levels will lead to more steps belief propagation for sharing local information globally. In our model, we adopt two message passing steps, as shown in Figure 2.

**Multi-loss training:** Since the goal of our model is to recognize group activities through global features and individual actions in that group, we adopt an alternative strategy for training the model. For the *k*th message passing step, we first remove the loss layers for actions and poses to learn parameters for group activity classification alone. In this phase, there is no back-propagation on action and pose classification. Since group activity heavily depends on an individual's activity, we then fix the softmax loss layer for scene classification and learn the model for actions and poses. The trained model is used for the next message passing step. Note that in each message passing step, we exploit the benefit of the neural network structure and jointly trained the whole network.

Learning semantic features for group activity: Traditional convolutional neural networks mainly focus on learning features for basic classification or localization tasks. However, in our proposed message passing CNN deep model, we not only learn features, but also learn semantic high-level features for better representing group activities and interactions within the group. We explore different layers' features for this deep model, and results show that these semantic features can be used for better scene understanding and classification.

**Implementation details:** Firstly, in practice, it is not guaranteed that every frame has the same number of detections. However, the structure of neural network should be fixed. To solve this problem, denoting  $M_{max}$  as the maximum number of people contained in one frame, we do a dummy-image padding when the number of people is less than  $M_{max}$ . Then we filter out these dummy data by de-activating neurons connected with them in related layers. Secondly, After the first message passing step, instead of directly feeding the raw scores into the next message passing step, we first normalize the pose and action scores for each person and scene scores for one frame by a softmax layer, converting to probabilities similar to belief propagation.

### 4 **Experiments**

Our models are implemented using the Caffe library [III] by defining two types of sparsely connected and weight shared inner product layers. One is from variable nodes to factor nodes, another is the reverse direction. We used TanH neurons as the non-linearity of these two layers. To examine the performance of our model, we test our model for scene classification on two datasets: (1) Collective Activity [II], (2) a nursing home dataset consisting of surveillance videos collected from a nursing home.

We trained an RBF kernel SVM on features extracted from the graphical model layer after each step of message passing model. These SVMs are used to predict scene labels for each frame, the standard task in these datasets.

#### 4.1 Collective Activity Dataset

The Collective Activity Dataset contains 44 video clips acquired using low resolution handheld cameras. Every person is assigned one of the following five action labels: crossing, waiting, queuing, walking and talking and one of the eight pose labels: right, front-right, front, front-left, left, back-left, back, back-right. Each frame is assigned one of the following five activities: crossing, waiting, queueing, walking, and talking. The activity category is attained by taking the majority of actions happening in one frame while ignoring the poses. We adopt the standard training test split used in [**L**].

In the Collective Activity dataset experiment, we further concatenate the global features for a scene with AC descriptors by HOG features [12]. We simply averaged AC descriptors features for all people and use this feature to serve as additional global information, namely this feature does not truly participated in the message passing process. This additional global information assists in classification with the limited amount of training data available for this dataset<sup>3</sup>.

We summarize the comparisons of activity classification accuracies of different methods in Table 1. The current best result using spatial information in graphical model is 79.1%, from Lan *et al.* [ $\square$ ], which adopted a latent max-margin method to learn graphical model with optimized structure. Our classification accuracies (the best is 80.6%) are competitive

<sup>&</sup>lt;sup>3</sup>Scene classification accuracy solely using AlexNet is 48%.

Z. DENG ET AL.: DEEP STRUCTURED MODELS FOR GROUP ACTIVITY RECOGNITION 9



Figure 4: Results visualization for our model. Green tags are ground truth, yellow tags are predicted labels. From left to right is without message passing, first step message passing and second step message passing

compared with the state-of-the-art methods. However, the benefits of the message passing are clear. Through each step of the message passing, the factor layer effectively captured dependencies between different variables and passing messages using factor neurons results in a gain in classification accuracy. Some visualization results are shown in Figure 4.

	1 Step MP	2 Steps MP	Latent Constituent [1]	75.1%
Pure Deep Learning (DL)	73.6%	78.4%	Contextual model [12]	79.1%
SVM+DL Feature	75.1%	80.6%	Our Best Result	80.6%

Table 1: Scene classification accuracy on the Collective Activity Dataset.

### 4.2 Nursing Home Dataset

This dataset consists 80 videos and is captured in a nursing home, including a variety of rooms such as dining rooms, corridors, etc. The 80 surveillance videos are recorded at 640 by 480 pixels at 24 frames per second, and contain a diverse set of actions and frequent cluttered scenes. This dataset contains typical actions include walking, standing, sitting, bending, squatting, and falling. For this dataset, the goal is to detect falling people, thus we assign each frame one of two activity categories: fall and non-fall. A frame is assigned "fall" if any person falls and "non-fall" otherwise. Note that many frames are challenging, and the falling person may be occluded by others in the scene. We adopted a standard 2/3 and 1/3 training test split. In order to remove redundancy, we sampled 1 out of every 10 frames for training and evaluation. Since this dataset has a large intra-class diversity within actions, we used the action primitive based detectors proposed in [LS] for more robust detection results.

Note that since this dataset has no pose attribute, we used the interaction between scene and actions instead to perform the two step message passing. For the SVM classifier, only deep learning features are used. We summarize the comparisons of activity classification accuracies of different methods in Table 2.

Ground Truth	Pure DL	SVM+DL Fea.	Detection	Pure DL	SVM+DL Fea.
1 Step MP	82.5%	82.3%	1 Step MP	74.4%	76.5%
2 Steps MP	84.1%	84.7%	2 Steps MP	75.6%	77.3%
<b>F</b> 11 <b>0 C</b> 1 <b>C C</b>			.1	II D	

Table 2: Classification accuracy on the Nursing Home Dataset

The scene classification accuracy on the Nursing Home dataset by using a baseline AlexNet model is 69%. The results on scene classification for each step also shows gains. In this dataset, accuracy on the second message passing gains an increase of around 1.5% for both pure deep learning or SVM prediction. We believe that this is due to the fact that the dataset only contains two scene labels, fall or non-fall, so scene variables are not as informative as scenes in the Collective Activity Dataset. Note that in both datasets, performance of scene classification plateaued after the second step message passing.

# 5 Conclusion

We have presented a deep learning model for group activity recognition which jointly captures the group activity, the individual person actions, and the interactions between them. We propose a way to combine graphical models with a deep network by mimicking the message passing process to do the inference mechanism. The model was successfully applied to real surveillance videos and the experiments showed the effectiveness of our approach in recognizing activities of a group of people.

## References

- [1] Ilya Sutskever Alex Krizhevsky and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems (NIPS)*, 2012.
- [2] Mohamed R. Amer, Dan Xie, Mingtian Zhao, Sinisa Todorovic, and Song-Chun Zhu. Cost-sensitive top-down / bottom-up inference for multiscale activity recognition. In European Conference on Computer Vision (ECCV), 2012.
- [3] Mohamed Rabie Amer, Peng Lei, and Sinisa Todorovic. Hirf: Hierarchical random field for collective activity recognition in videos. In *European Conference on Computer Vision (ECCV)*, pages 572–585, 2014.
- [4] Borislav Antic and BjÄűrn Ommer. Learning latent constituents for recognition of group activities in video. In *European Conference on Computer Vision (ECCV)*, 2014.
- [5] W. Brendel and S. Todorovic. Learning spatiotemporal graphs of human activities. In *International Conference on Computer Vision (ICCV)*, pages 778–785, 2011.
- [6] W. Choi and S. Savarese. A unified framework for multi-target tracking and collective activity recognition. In *European Conference on Computer Vision (ECCV)*, 2012.

#### Z. DENG ET AL.: DEEP STRUCTURED MODELS FOR GROUP ACTIVITY RECOGNITION 11

- [7] Wongun Choi, Khuram Shahid, and Silvio Savarese. What are they doing?: Collective activity classification using spatio-temporal relationship among people. In *International Conference on Computer Vision Workshops on Visual Surveillance*, pages 1282– 1289. IEEE, 2009.
- [8] Jia Deng, Nan Ding, Yangqing Jia, Andrea Frome, Kevin Murphy, Samy Bengio, Yuan Li, Hartmut Neven, and Hartwig Adam. Large-scale object classification using label relation graphs. In *European Conference on Computer Vision (ECCV)*, pages 48–64. Springer, 2014.
- [9] Abhinav Gupta, Praveen Srinivasan, Jianbo Shi, and Larry S. Davis. Understanding videos, constructing plots: Learning a visually grounded storyline model from annotated videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [10] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [11] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and Li Fei-Fei. Largescale video classification with convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR)*, pages 1725–1732, 2014.
- [12] Tian Lan, Wang Yang, Yang Weilong, and Greg Mori. Beyond actions: Discriminative models for contextual group activities. In *Advances in Neural Information Processing Systems (NIPS)*, 2010.
- [13] Tian Lan, Leonid Sigal, and Greg Mori. Social roles in hierarchical models for human activity recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [14] Tian Lan, Yang Wang, Weilong Yang, Stephen Robinovitch, and Greg Mori. Discriminative latent models for recognizing contextual group activities. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 34(8):1549–1562, 2012.
- [15] Tian Lan, Chen Lei, Deng Zhiwei, Guang-Tong Zhou, and Greg Mori. Learning action primitives for multi-level video event understanding. In *International Workshop on Visual Surveillance and Re-Identification (at ECCV)*, 2014.
- [16] R. Poppe. A survey on vision-based human action recognition. IVC, 28:976–990, 2010.
- [17] Vignesh Ramanathan, Bangpeng Yao, and Li Fei-Fei. Social role discovery in human events. In *Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [18] M. S. Ryoo and J. K. Aggarwal. Stochastic representation and recognition of high-level group activities. *International Journal of Computer Vision (IJCV)*, 93(2):183–200, 2011.
- [19] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems* (*NIPS*), pages 568–576. Curran Associates, Inc., 2014.

#### 12 Z. DENG ET AL.: DEEP STRUCTURED MODELS FOR GROUP ACTIVITY RECOGNITION

- [20] J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [21] K.N. Tran, A. Gala, I.A. Kakadiaris, and S.K. Shah. Activity analysis in crowded environments using social cues for group discovery and human interaction modeling. *Pattern Recognition Letters*, 44:49–57, 2014.
- [22] Heng Wang and C. Schmid. Action recognition with improved trajectories. In *International Conference on Computer Vision (ICCV)*, pages 3551–3558, 2013.
- [23] D. Weinland, R. Ronfard, and E. Boyer. A survey of vision-based methods for action representation, segmentation and recognition. In *Computer Vision and Image Understanding (CVIU)*, 2010.
- [24] Y. Zhu, N. Nayak, and A. Roy-Chowdhury. Context-aware modeling and recognition of activities in video. In *Computer Vision and Pattern Recognition (CVPR)*, 2013.