

Understanding Social Relationship with Person-pair Relations

Sarit Kraus

Department of Computer Science, Bar-Ilan University, Israel
pcchair@ijcai19.org

Abstract

Social relationships understanding is to infer the social relations among people from images and videos, which has attracted increasing attention in computer vision recently. A great progress has been made since the rise of deep learning. However, they mostly focus on the facial attributes or contextual object cues without taking into account the interaction among person pairs. Motivated by scene graph generation, we carefully analyzed the datasets and found the social relations in a still image always have high semantic relevance. For instance, if two person pair in an image are *Friends*, then the third one is always friends or at least other intimate relations but not *No Relation*. Therefore, to capture this interaction cues, we propose a novel end-to-end trainable Person-Pair Relation Network (PRN) using standard RNNs, a graph inference network that learns iteratively to improve its predictions via message passing among person pair nodes. Extensive experiments on PISC and PIPA-Relation show the superiority of our method over previous methods.

1 Introduction

Social relationships are closely related to our daily life [Barr *et al.*, 2014]. After understanding the social relationship between the person pair, we can easily explain their behavior. For machines, only when they fully understand the social relationships, can they further understand and infer the human behavior in our social life, so as to make a better response. In addition, we often leave traces that capture social relationships in many medias and we not only want the machines to be proficient at their task, but also enable them to blend in and act appropriately in different situations [Sun *et al.*, 2017]. In short, social relationship detection task is very significant in many ways. In our work, we aim to address the social relationship detection task for every picture where each picture represents a scene.

However, to solve the social relationship detection task is not so simple. For a giving picture, detecting the social relationships of all the person pair is a difficult task. The models

need to be adapted to different scenes and context information to make right judgments. [Sun *et al.*, 2017] use the information of head region, body region and human attributes to predict the person-pair’s social relationship separately. [Li *et al.*, 2017] make use of the pair of people in question and region proposals and allocate attention to each region to detect the social relationship of each person pair. [Wang *et al.*, 2018] takes advantage of the message propagation between person pair social relationship and the object semantic regions to solve the problem. The biggest problem of these models is that they all only detect one relationship per step which will cause that different social relationships in the same scene cannot interact with each other. Social relationships in the same scene are strongly linked but the previous models have ignored this important information.

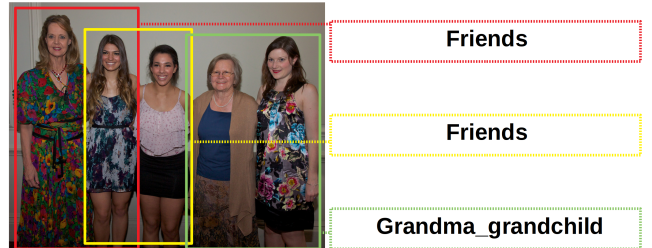


Figure 1: An example image from PIPA dataset. There are 3 person-pair social relationships in the picture: friends, friends and grandma_grandchild.

As one example in PIPA (Figure 2) where the picture denotes a scene, there are 2 "friends" relationships and 1 "grandma_grandchild" relationship in the scene. If we know only that two of these relationships are "friends," we can infer that it is most likely a group of friends, and thus we infer that the unknown relationship is "friends" if they're similar in age. In other words, we can easily use the interaction of different relationships in the same scene to detect each social relationship. Therefore we address the issue with focusing more on the interaction between every social relationships in one scene to improve the social relationship detection performance.

However, the biggest issue of the task is that it is not as simple as people directly to judge the result. We need to design the mechanism of the interaction between social relationships and effectively model the mechanism. The other

issue is how to use less effective information to model the interaction mechanism and get the better result. For example, whether to introduce the information of the objects in scene has also become an important consideration for us.

To address this problem, we propose a novel end-to-end trainable Person-Pair Relation Network (PRN) which makes good use of the information of social relationship interaction in the same scene. PRN consists of 3 parts: Feature Extraction Module, Message Pooling Module and Message Passing Module. In the Feature Extraction Module, we use a Resnet to extract features of each person and another Resnet to extract features of the person pairs in the scene. In this module, the position of each person is also taken into account. In the Message Pooling Module, we use a pooling mechanism to make all the social relationships interact with each other well and the output of this module will be considered as the inputs of the GRU which is in the next module. In the Message Passing Module, we use the RNNs contained GRU to make the social relationship message passing proceed iteratively. As the model proceeding, this module and the Message Pooling Module will have an interaction and make the overall interaction mechanism better. We use the hidden state of the last GRU as the social relationship detection results of the scene.

We evaluate our model in classic datasets: PIPA-Relation dataset and PISC dataset. The experiment results verify the superiority of our model over previous methods. Specifically, our contributions are as follows:

- To our best knowledge, it is the first attempt to introduce the interaction of social relationships in the same scene on the social relationship detection task;
- We design a novel interaction mechanism to model the interaction of social relationships in the same scene which gets the best result in this task;
- We analyze and verify that the role played by objects in the scene is not very big which is a very novel idea in the social relationship detection task.

2 Related Work

2.1 Social Relationship Understanding

The foundation of social network is the social relationships understanding, an important multidisciplinary problem that has attracted increasing attention in computer vision recently. A much number of studies that aim to infer social relationships from images [Wang *et al.*, 2015; Li *et al.*, 2017; Wang *et al.*, 2018; 2010; Zhang *et al.*, 2015b] and videos [Ding and Yilmaz, 2010; Ramanathan *et al.*, 2013; Vinciarelli *et al.*, 2009] have been made since the rise of deep learning. For instance, motivated by psychological studies, [Zhang *et al.*, 2015b] and [Dibeklioglu *et al.*, 2013] exploit social relationships based on facial attributes such as expression and head pose, and affective behaviour analysis. Besides, [Li *et al.*, 2017] and [Wang *et al.*, 2018] discover that contextual cues around people play a significant role in social relationship inferring. Concretely, [Li *et al.*, 2017] proposed a dual-glance model for social relationship, where the first glance makes a coarse relationship prediction for a given person pair and

then the second one refines the prediction by using the objects around the pair. [Wang *et al.*, 2018] constructed a semantic-aware knowledge graph and employed Gated Graph Neural Network (GGNN) [Li *et al.*, 2015] to integrate the graph into the Graph Reasoning Model (GRM), a graph reasoning network where a proper message propagation and graph attention mechanism are introduced to explore the interaction between person pair and the contextual objects.

Unlike the aforementioned works which mainly focus on facial attributes or contextual object cues, we detailly studied the two classic datasets PISC [Li *et al.*, 2017] and PIPA-relation [Sun *et al.*, 2017] and found the social relations in a still image have high semantic relevance. Based on this discovery, we designed a novel end-to-end trainable Person-Pair Relation Network (PRN), a graph inference network to capture this semantic relevance cues via message passing among person pair nodes.

2.2 Message Passing

Graph inference is a kind of form of message passing and Conditional Random Fields (CRF) have been used extensively in this field. Johnson *et al.* used CRF to infer scene graph grounding distributions for image retrieval [Johnson *et al.*, 2015]. Yatskar *et al.* use a deep CRF model to propose situation-driven object and action prediction [Yatskar *et al.*, 2016]. Danfei Xu *et al.* use the GRU-RNNs to solve the scene graph generation problem iteratively [Xu *et al.*, 2017]. Our work is related to Graph-LSTM [Liang *et al.*, 2016] and the work of Danfei Xu *et al.* [Xu *et al.*, 2017] which formulate the message passing problem using RNN models. Danfei Xu *et al.* [Xu *et al.*, 2017] design primal graph and dual graph in their model while we just simplify the model and just use one graph to make social relationship messages to pool and achieve a better result. As what Danfei Xu *et al.* [Xu *et al.*, 2017] have done, our model iteratively refines the social relationship predictions through relationship message passing in the scene, whereas the Structural RNN model only makes one-time predictions along the temporal dimension, and thus cannot refine its past predictions [Xu *et al.*, 2017].

3 PRN model

This part will be written by **liangjinrui** and **chenhaicheng**.
[model figure]

[Introduce the total model]

3.1 Social Relationship Understanding Model

3.2 Imposing Rules

3.3 Integrating by Integer Linear Programming

3.4 Optimization

Given the predicted score $s^{I,k} \in R^{|C|}$ for the k -th person pair in image I , we use *softmax* to get its corresponding probability $p^{I,k} \in R^{|C|}$

$$p_i^{I,k} = \frac{\exp s_i^{I,k}}{\sum_{j=1}^{|C|} \exp s_j^{I,k}}, i = 1, 2, \dots, |C| \quad (1)$$

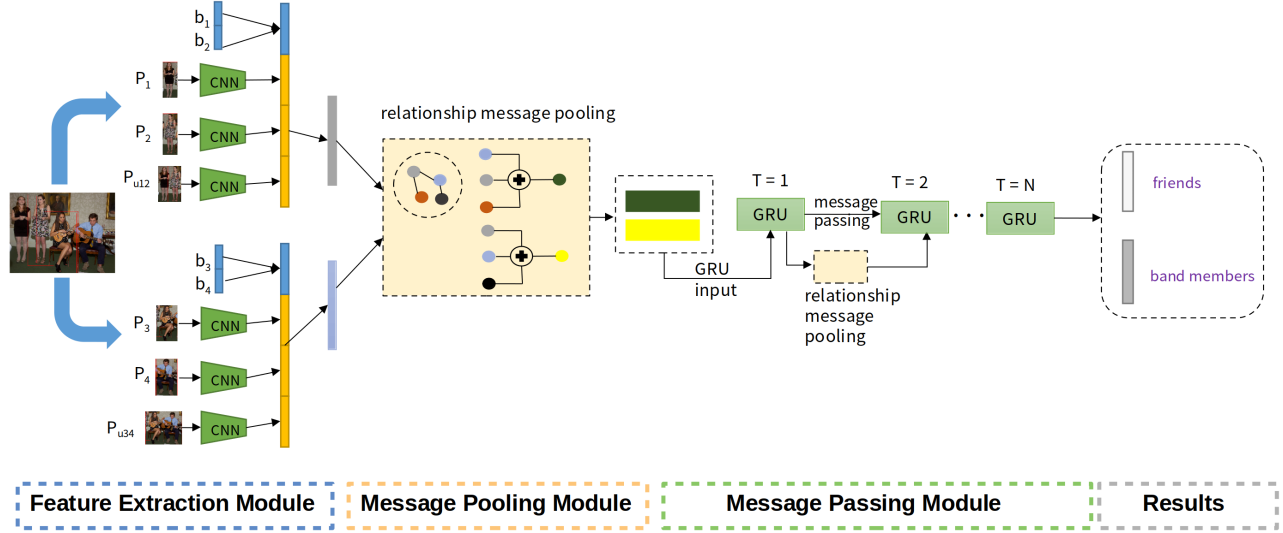


Figure 2: An illustration of our PRN. The model first extracts the features of each person and person pairs in the Feature Extraction Module. P_i denotes the i th person and P_{uij} denotes the union region of i th person and j th person. b_i denotes the position of the i th person. In the Message Pooling Module, a message pooling function computes relationship messages that are passed to the node GRU in the next iteration from the hidden states. The \oplus symbol denotes a learnt weighted sum. Then in the Message Passing Module, we iteratively updates the hidden states of the GRUs. We use the hidden states of the GRUs at the last iteration step, to predict the social relationships in the scene.

where \mathcal{C} donates the classes set of social relationship and $|\mathcal{C}|$ is its size. The loss function is expressed as

$$\mathcal{L} = -\frac{1}{\sum_{I \in \mathcal{I}} N(I)} \sum_{I \in \mathcal{I}} \sum_{k=1}^{N(I)} \sum_{i=1}^{|\mathcal{C}|} L(y_i^{I,k}, p_i^{I,k}) \quad (2)$$

where $N(I)$ returns the number of person pair in image I , $L(\cdot)$ is the cross entropy loss function, \mathcal{I} is the image set.

4 Experiments

This part will be written by **chenhaicheng**.

4.1 Experiment Setting

Datasets. In this work, two datasets were used to evaluate our proposed method and other existing ones. The first one is the large-scale People in Social Context (PISC) [Li *et al.*, 2017] with 22,670 images and contains two-level recognition tasks: **3 Coarse-level relationship**, namely *No Relation*, *Intimate Relation*, *None-Intimate Relation* and **6 Fine-level relationship**, i.e., *Friend*, *Family*, *Couple*, *Professional*, *Commerical*, *No Relation*. The second one is the People in Photo Album Relation (PIPA-Relation) [Sun *et al.*, 2017], an extension version of People in Photo Album (PIPA) [Zhang *et al.*, 2015a] with 37107 images. It also annotates 26,915 person pairs on two-level recognition tasks: **5 Social Domains** and **16 Social Relations** based on these domains. The train/val/test in PISC are 13,142/4,000/4,000 images with 14,536/25,636/15,497 person pairs on coarse level relationship, and 16,828/500/1,250 images with 55,400/1,505/3,691 person pairs on fine level relationship, respectively. In PIPA-Relation, we follow [Wang *et al.*, 2018] and focus on recognizing its 16 relationships in the experiment. The train/val/test in it are 13,729/709/5,106 person pairs.

Implementation Details. During our work, We adopt the same strategy as previous works including [Li *et al.*, 2017] and [Wang *et al.*, 2018]. First, we fine-tune the ResNet-101 model [He *et al.*, 2016], and we set the a lower learning rate as 0.0001. For the message passing propagation model, the dimension of hidden size is set as 512. The iteration time T is set as 4 and learning rate as 0.0001. Similar to [Wang *et al.*, 2018], we the fine-tuning model utilized SGD, and the message passing module is trained with ADAM.

4.2 Datasets Analysis

In this subsection, we made an analysis on PISC and PIPA-Relation. Here we only used its train set and test set for statistic. As shown in Figure 3, we first calculated the social relation categories of each image, and found almost all images have only one social relation category, and the other half have two categories. For example, on PISC, approximately 79.944% of images have only one coarse category, while 20.0327% images have two coarse one.

4.3 Comparisons with State-of-the-Art Methods

We compare our proposed model with existing state-of-the-art methods on both PISC and PIPA-Relation datasets. Formally, the compared methods are as followed:

Performance on the PISC dataset

UnionCNN Following [Lu *et al.*, 2016], it generates a single CNN model to predicate relations. In this task, we also feeds the union region of person pair to s single CNN for classification.

Pair CNN [Li *et al.*, 2017] consists of two equivalent CNNs with shared weights to extracted features for image for two individuals.

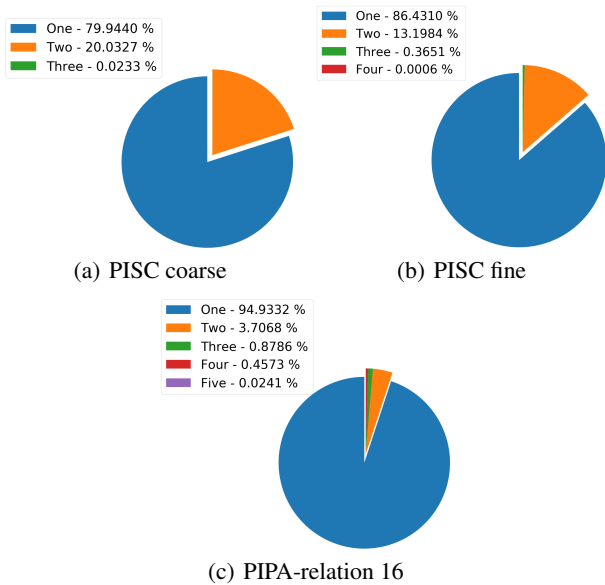


Figure 3: Social relation categories per image on PISC and PIPA-relation.

Pair CNN + BBox + Union[Li *et al.*, 2017] incorporates spatial location information of two bounding box that based the previous pair CNN and Union CNN.

Dual-glance[Li *et al.*, 2017] implements coarse and fine prediction which includes three and six relationships. Dual-glance employing pair CNN + BBox + BBox + Union and utilized surrounding region proposal to refine the prediction.

GRM[Wang *et al.*, 2018] propose a graph reasoning model that unifies the frequency of co-concurrences of each relationship-object pair to facilitate social relation. Similar to the model of GRM, we also adopt the per-class recall and mean average precision(mAP) to evaluate our model. The experiments data are reported in Table 1. First, both Pair CNN + BBox + Union, Pair CNN + BBox + Global, Dual-glance are incur extra Faster-RCNN[Ren *et al.*, 2015] to extract the local contextual cues(object proposal). GRM utilized the object proposal to construct a semantic-aware knowledge graph to reason about the social relationship. It is notable that both of them incur extra detection annotations that contains noises. dasdasdasd

Performance on the PIPA-Relation dataset

On this dataset, we also compare our proposed model with the existing Two stream CNN[Sun *et al.*, 2017],Dual-glance[Li *et al.*, 2017] and GRM[Wang *et al.*, 2018].Two stream CNN takes

4.4 Experiment Results

4.5 Experiment Analysis

4.6 Ablation Study

4.7 Case Study

Two examples in Figure 4 are shown to illustrate the capacity of our PRN to infer social relationships. We compare two pre-

Figure 4: Comparison among social relationships with different sources, the green rectangle from the dataset, the yellow from our PRN and the blue from GRM. The person with box in various colors on left original image corresponds to the same color node in the right.

dicted social relationships, one from our PRN and the another from GRM [Wang *et al.*, 2018].

5 Conclusion

This part will be written by **liangjinrui**.

Table 1: Recall-per-class and mean average precision (mAP) evaluating our PRN model and previous methods on PISC (in %).

Methods	Coarse relationships				Fine relationships						
	Intimate	Non-Intimate	No Relation	mAP	Friends	Family	Couple	Professional	Commerical	No Relation	mAP
Union CNN [Lu <i>et al.</i> , 2016]	72.1	81.8	19.2	58.4	29.9	58.5	70.7	55.4	43.0	19.6	43.5
Pair CNN [Li <i>et al.</i> , 2017]	70.3	80.5	38.8	65.1	30.2	59.1	69.4	57.5	41.9	34.2	48.2
Pair CNN + BBox + Union [Li <i>et al.</i> , 2017]	71.1	81.2	57.9	72.2	32.5	62.1	73.9	61.4	46.0	52.1	56.9
Pair CNN + BBox + Global [Li <i>et al.</i> , 2017]	70.5	80.0	53.7	70.5	32.2	61.7	72.6	60.8	44.3	51.0	54.6
Dual-glance [Li <i>et al.</i> , 2017]	73.1	84.2	59.6	79.7	35.4	68.1	76.3	70.3	57.6	60.9	63.2
GRM [Wang <i>et al.</i> , 2018]	81.7	73.4	65.5	82.8	59.6	64.4	58.6	76.6	39.5	67.7	68.7
Ours											

Table 2: Accuracy (in %) evaluating our PRN model and previous methods on PIPA-relation.

Methods	accuracy
Two stream CNN [Zhang <i>et al.</i> , 2015a]	57.2
Dual-Glance [Li <i>et al.</i> , 2017]	59.6
GRM [Wang <i>et al.</i> , 2018]	62.3
Ours	

References

- [Barr *et al.*, 2014] Jeremiah R. Barr, Leonardo A. Cament, Kevin W. Bowyer, and Patrick J. Flynn. Active clustering with ensembles for social structure extraction. In *IEEE Winter Conference on Applications of Computer Vision, Steamboat Springs, CO, USA, March 24-26, 2014*, pages 969–976, 2014.
- [Dibeklioglu *et al.*, 2013] Hamdi Dibeklioglu, Albert Ali Salah, and Theo Gevers. Like father, like son: Facial expression dynamics for kinship verification. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 1497–1504, 2013.
- [Ding and Yilmaz, 2010] Lei Ding and Alper Yilmaz. Learning relations among movie characters: A social network perspective. In *Computer Vision - ECCV 2010, 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part IV*, pages 410–423, 2010.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016.
- [Johnson *et al.*, 2015] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Image retrieval using scene graphs. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3668–3678, 2015.
- [Li *et al.*, 2015] Li-Jia Li, David A. Shamma, Xiangnan Kong, Sina Jafarpour, Roelof van Zwol, and Xuanhui Wang. Celebritynet: A social network constructed from large-scale online celebrity images. *TOMCCAP*, 12(1):3:1–3:22, 2015.
- [Li *et al.*, 2017] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S. Kankanhalli. Dual-glance model for deciphering social relationships. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2669–2678, 2017.
- [Liang *et al.*, 2016] Xiaodan Liang, Xiaohui Shen, Jiashi Feng, Liang Lin, and Shuicheng Yan. Semantic object parsing with graph LSTM. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 125–143, 2016.
- [Lu *et al.*, 2016] Cewu Lu, Ranjay Krishna, Michael S. Bernstein, and Fei-Fei Li. Visual relationship detection with language priors. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 852–869, 2016.
- [Ramanathan *et al.*, 2013] Vignesh Ramanathan, Bangpeng Yao, and Fei-Fei Li. Social role discovery in human events. In *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, pages 2475–2482, 2013.
- [Ren *et al.*, 2015] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99, 2015.
- [Sun *et al.*, 2017] Qianru Sun, Bernt Schiele, and Mario Fritz. A domain based approach to social relation recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 435–444, 2017.

- [Vinciarelli *et al.*, 2009] Alessandro Vinciarelli, Maja Pantic, and Hervé Bourlard. Social signal processing: Survey of an emerging domain. *Image Vision Comput.*, 27(12):1743–1759, 2009.
- [Wang *et al.*, 2010] Gang Wang, Andrew C. Gallagher, Jiebo Luo, and David A. Forsyth. Seeing people in social context: Recognizing people and social relationships. In *Computer Vision - ECCV 2010 - 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part V*, pages 169–182, 2010.
- [Wang *et al.*, 2015] Quan Wang, Bin Wang, and Li Guo. Knowledge base completion using embeddings and rules. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 1859–1866, 2015.
- [Wang *et al.*, 2018] Zhouxia Wang, Tianshui Chen, Jimmy S. J. Ren, Weihao Yu, Hui Cheng, and Liang Lin. Deep reasoning with knowledge graph for social relationship understanding. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 1021–1028, 2018.
- [Xu *et al.*, 2017] Danfei Xu, Yuke Zhu, Christopher B. Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 3097–3106, 2017.
- [Yatskar *et al.*, 2016] Mark Yatskar, Luke S. Zettlemoyer, and Ali Farhadi. Situation recognition: Visual semantic role labeling for image understanding. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 5534–5542, 2016.
- [Zhang *et al.*, 2015a] Ning Zhang, Manohar Paluri, Yaniv Taigman, Rob Fergus, and Lubomir D. Bourdev. Beyond frontal faces: Improving person recognition using multiple cues. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 4804–4813, 2015.
- [Zhang *et al.*, 2015b] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Learning social relation traits from face images. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 3631–3639, 2015.