

# From Facial Expression Recognition to Interpersonal Relation Prediction

Zhanpeng Zhang · Ping Luo · Chen Change Loy · Xiaoou Tang

Received: date / Accepted: date

**Abstract** Interpersonal relation defines the association, *e.g.*, warm, friendliness, and dominance, between two or more people. We investigate if such fine-grained and high-level relation traits can be characterized and quantified from face images in the wild. We address this challenging problem by first studying a deep network architecture for robust recognition of facial expressions. Unlike existing models that typically learn from facial expression labels alone, we devise an effective multitask network that is capable of learning from rich auxiliary attributes such as gender, age, and head pose, beyond just facial expression data. While conventional supervised training requires datasets with complete labels (*e.g.*, all samples must be labeled with gender, age, and expression), we show that this requirement can be relaxed via a novel attribute propagation method. The approach further allows us to leverage the inherent correspondences between heterogeneous attribute sources despite the disparate distributions of different datasets. With the network we demonstrate state-of-the-art results on existing facial expression recognition benchmarks. To predict interpersonal relation, we use the expression recognition network

as branches for a Siamese model. Extensive experiments show that our model is capable of mining mutual context of faces for accurate fine-grained interpersonal prediction.

**Keywords** Facial Expression Recognition · Interpersonal Relation · Deep Convolutional Network

## 1 Introduction

Facial expression recognition is an actively researched topic in computer vision [70]. Existing pipelines typically recognize single-person expressions and assign them into discrete prototypical classes, namely anger, disgust, fear, happy, sad, surprise, and neutral. Inspired by extensive psychological studies [19,21,23,34], in this work we wish to investigate the interesting problem of characterizing and quantifying interpersonal relation traits from human face images beyond just expressions.

Interpersonal relation manifests when one establish, reciprocate, or deepen relationships with one another. The recognition task goes beyond facial expression recognition that analyzes facial motions and facial feature changes of a single subject. It aims for a higher-level interpretation of fine-grained and high-level interpersonal relation traits, such as friendliness, warm, and dominance for faces that co-exist in an image. Effectively exploiting such relational cues can provide rich social facts. An example is shown in Fig. 1. Such a capability promises a wide spectrum of applications. For instance, automatic interpersonal relation inference allows for relation mining from image collection in social networks, personal albums, and films. Face-based relational cues can also be combined with other visual cues

---

Zhanpeng Zhang  
SenseTime Group Limited  
E-mail: zhangzhanpeng@sensetime.com

Ping Luo  
Department of Information Engineering, The Chinese University of Hong Kong  
E-mail: pluo@ie.cuhk.edu.hk

Chen Change Loy (Corresponding Author)  
Department of Information Engineering, The Chinese University of Hong Kong  
E-mail: ccloy@ie.cuhk.edu.hk

Xiaoou Tang  
Department of Information Engineering, The Chinese University of Hong Kong  
E-mail: xtang@ie.cuhk.edu.hk



**Fig. 1:** The image is given a caption ‘German Chancellor Angela Merkel and U.S. President Barack Obama inspect a military honor guard in Baden-Baden on April 3.’ (source: [www.rferl.org](http://www.rferl.org)). When we examine the face images jointly, we could observe far more rich social facts that are different from that expressed in the text.

such as body postures [7] to achieve an even richer modeling and prediction of relations<sup>1</sup>.

Profiling unscripted interpersonal relation from face images is non-trivial. Among the most significant challenges are:

1. Most existing face analysis models only consider a single subject. No existing methods attempt to consider pairwise faces jointly.
2. Relations are governed by a number of high-level facial factors [19, 21, 23]. Thus we need a rich face representation that captures various attributes such as expression, gender, age, and head pose;
3. No single dataset is presently available to encompass all the required facial attribute annotations for learning such a rich representation. In particular, some datasets only contain face expression labels, while other datasets may only be annotated with the gender label. Moreover, these datasets are collected from different environments and exhibit vastly different statistical distributions. Model training on such heterogeneous data remains an open problem.

We address the first problem through formulating a novel deep convolutional network with a Siamese-like architecture [3]. The architecture consists of two convolutional network branches with shared parameters. Each branch is dedicated to one of the faces that co-exist in an image. Outputs of these two branches are fused to allow joint relation reasoning from pairwise faces, where each face serves as the mutual context to the other.

To address the second challenge, we formulate the convolutional network branches in a multitask framework such that it is capable of learning rich face representation from auxiliary attributes such as head pose, gender, and age, apart from just facial expressions. To facilitate the multitask learning, we gather various existing face expression and attribute datasets and additionally label a new large-scale face

**Expression in-the-Wild (ExpW)** dataset, which is formed by over 90,000 web images.

To mitigate the third issue of learning from heterogeneous datasets, we devise a new attribute propagation approach that is capable of dealing with missing attribute labels from different datasets, and yet bridging the gap of heterogeneous datasets. In particular, during the training process, our network dynamically infers missing attribute labels of a sample using Markov Random Field (MRF), conditioned on appearance similarity of that sample with other annotated samples. We will show that the attribute propagation approach allows our network to learn effectively from heterogeneous datasets with different annotations and statistical distributions.

The contributions of this study include:

1. We make the first attempt to investigate face-driven fine-grained interpersonal relation prediction, of which the relation traits are defined based on psychological study [33]. We carefully investigate the detectability and quantification of such traits from face image pairs.
2. We formulate a new deep architecture for learning face representation driven by multiple tasks, *e.g.* pose, expression, and age. Specifically, we introduce a new attribute propagation approach to bridge the gap from heterogeneous sources with potentially missing target attribute labels. We show that this network leads to new state-of-the-art results on widely-used facial expression benchmarks. It also establishes a solid foundation for us to recognize interpersonal relations.
3. We construct a new interpersonal relation dataset labeled with pairwise relation traits supported by psychological studies [33, 34]. In addition, we also introduce a large-scale facial expression in-the-wild dataset<sup>2</sup>.

In comparison to our earlier version of this work [93], we present a more principle and unified way of addressing the heterogeneous data problem using the MRF-based attribute propagation approach. This is in contrast to the deep bridging layer proposed in our previous work [93], which requires external facial alignment step to extract local part appearances for establishing cross-dataset association. In addition, we study more closely on the facial expression recognition problem, which is crucial for accurate interpersonal relation identification. Specifically, we present a new large-scale dataset and conduct extensive experiments against state-of-the-art expression recognition methods. Apart from the methodology, the paper was also substantially improved by providing more technical details and more extensive experimental evaluations.

<sup>1</sup> Despite we did not study the integration of face and body cues, if body posture and hand gesture information are available, they can be naturally used as additional input channels for our deep models.

<sup>2</sup> Both ExpW and relation datasets are available at <http://mmlab.ie.cuhk.edu.hk/projects/socialrelation/index.html>

## 2 Related Work

Understanding interpersonal relation can be regarded as a subfield under *social signal processing* [8,56,60,75,76], an important multidisciplinary problem that has attracted a surge of interest from computer vision community. Social signal processing mainly involves facial expression recognition [96,70,45,64,80,16,46,44,31,52,16,97]. We provide a concise account as follows.

**Facial expression recognition.** A facial expression recognition algorithm usually consists of face representation extraction and classifier construction. Depending on the adopted face representation, existing algorithms can be broadly categorized into two groups: facial action based methods and appearance-based approaches.

Facial action based methods usually exploit the face geometrical information or face action units driven representation for facial expression classification. For example, Tiam *et al.* [71] use the positions of facial landmarks for facial action recognition and then perform expression analysis. Ruiz *et al.* [64] combine the tasks of facial action detection and expression recognition to leverage their coherence. Liu *et al.* [43] construct a deep network to learn a middle representation known as Micro-Action-Pattern (MAP) representation, so as to bridge the semantic gap between low-level features and high-level expression concepts. Liu *et al.* [44] adapt 3D Convolutional Neural Network (CNN) to detect specific facial action parts to obtain discriminative part-based representation.

Appearance-based methods extract features from face patches or the whole face region. A variety of hand-crafted features have been employed, such as LBP [74,96], HOG [9], and SIFT [25] features. Recently, a number of methods [31,32,46,52,54,88,97] attempt to learn facial features directly from raw pixels by deep learning. Unlike methods based on hand-crafted features, a deep learning framework allows end-to-end optimization of feature extraction, selection, and expression recognition. Liu *et al.* [46] show the effectiveness of Boosted Deep Belief Network (BDBN) for end-to-end feature extraction and selection. More recent studies [97] adopt CNN architectures that permit feature extraction and recognition in an end-to-end framework. For instance, Yu *et al.* [88] employed an ensemble of multiple deep CNNs. Mollahosseini *et al.* [52] used three inception structures [69] in convolution for facial expression recognition. The Peak-Piloted Deep Network (PPDN) [97] is introduced to implicitly learn the evolution from non-peak to peak expressions. We introduce readers to a recent survey [89] focusing on deep learning-based facial behavior analysis.

Our approach is regarded as an appearance-based approach, but differs significantly from the aforementioned studies in that most existing approaches are based on single person, therefore, cannot be directly employed for interper-

sonal relation inference. In addition, these studies mostly focus on recognizing prototypical expressions. Interpersonal relation is far more complex involving many factors such as age and gender. Thus we need to consider more attributes jointly in our problem.

**Human interaction and group behavior analysis.** There exists a number of studies that analyze human interaction and group behavior from images and videos [13,14,17,62,63,79,18]. Many of these studies focus on the coarser level of interpersonal connection other than the one defined by Kiesler in the interpersonal circle [33]. For instance, Ding and Yilmaz [13] and Ricci *et al.* [63] only identify the social group (or jointly for estimate head and body orientations) without inferring the relation between individuals. Fathi *et al.* [17] only detect three social interaction classes, *i.e.*, ‘dialogue, monologue and discussion’. Wang *et al.* [77] define social relation by several social roles, such as ‘father-child’ and ‘husband-wife’. Chakraborty *et al.* [5] classify photos into classes such as ‘couple, family, group, or crowd’. Other related problems also include image communicative intents prediction [30] and social role inference [39], usually applied on news and talks shows [61], or meetings to infer dominance [27].

In comparison to the aforementioned studies [13,17], our work aims to recognize fine-grained and high-level interpersonal relation traits [33], rather than identify social group and roles. In addition, many of these studies did not use face images directly, but visual concepts [14] discovered by detectors or people spatial proximity in 2D or 3D spaces [6]. All these information sources are valuable for learning human interactions but we believe that face still serves a primary role in defining fine-grained and high-level interpersonal relation since face can reveal much richer information such as expression, age, and gender.

Other group behavior studies [10,24,28,35] mainly recognize action-oriented behaviors such as hugging, handshaking or walking, but not face-based interpersonal relations. Often, group spatial configuration and actions are exploited for recognition. Our study differs in that we aim at recognizing abstract relation traits from faces.

**Deep learning.** Deep learning has achieved remarkable success in many tasks of face analysis, *e.g.* face detection [84,85,42,82,55], face parsing [50,47], face landmark detection [92,99,73], face attribute recognition [48,78,26], face recognition [66,58,68], and face clustering [94]. However, deep learning has not yet been adopted for face-driven interpersonal relation mining that requires joint reasoning from multiple persons. In this work, we propose a deep model to capture complex facial attributes from heterogeneous datasets, and joint learning from face pairs. Although there are several algorithms [2,87,83] that perform training on heterogeneous datasets, most of these studies assume fixed image features and exploit the label

**Table 1:** A comparison of popular facial expression datasets and the proposed ExpW dataset.

Datasets	Quantity	Environment	Expression	Data format
JAFFE [51]	213 images from 10 subjects	lab	posed	256×256 gray scale image
MMI [57]	238 sequences from 28 subjects	lab	posed	720×576 RGB frames
Oulu-CASIA [95]	480 sequences from 80 subjects	lab	posed	320×240 RGB frames
CK+ [49]	593 sequences from 123 subjects	lab	posed	640×490 or 640×480 gray scale frames
FER [20]	35,587 images	wild	natural	48×48 gray scale image
SFEW [12]	1,635 images	wild	natural	720×576 RGB images
ExpW	91,793 images	wild	natural	Original web images

correlation for missing label propagation. Lee [40] proposes a deep learning algorithm that employs pseudo label to utilize the unlabeled data. But the pseudo label is simply generated by a pre-trained network using labeled data, thus the potential correlation between the labeled and unlabeled data is ignored. Our network also differs from the multitask network in [92], which assumes complete labels from all attributes and homogeneous data sources.

### 3 Face Expression and Interpersonal Relation Datasets

Before we describe our approach, we introduce two new datasets collected in this study.

#### 3.1 Face Expression Dataset

Research in face perception and emotion typically requires very large annotated datasets of images of facial expressions. There are a number of facial expression datasets, *e.g.*, CK+ [49], JAFFE [51], Oulu-CASIA [95], MMI [57], FER [20], SFEW [12]. A summary is provided in Table 1. These datasets are either collected in controlled environments, or the quantity is insufficient to train a robust deep network. An automatic method for expression dataset construction is proposed in [16]. This method is useful to collect large-scale dataset. Nonetheless, it relies on accurate facial landmark detection and thus may limit face variations in the collected data.

To this end, we built a new database named as Expression in-the-Wild (ExpW) dataset that contains 91,793 faces manually labeled with expressions. The quantity of images in ExpW is larger and the face variations are more diverse than many existing databases, as summarized in Table 1. Figure 2 shows some example images of ExpW.

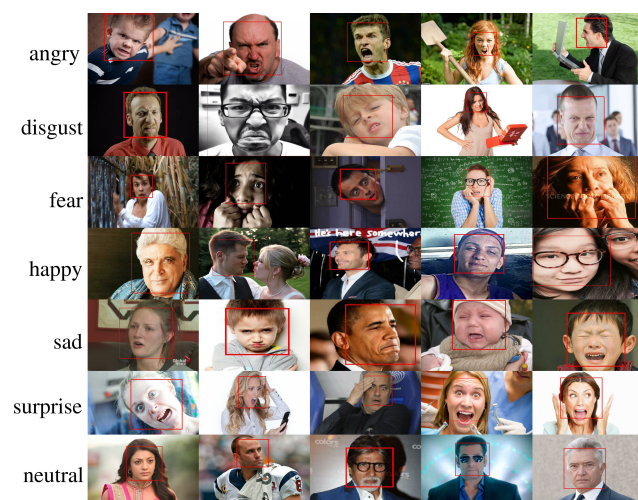
We collected ExpW dataset in the following way. Firstly, we prepared a list of emotion-related keywords such as “excited”, “afraid” and “panic”. Then we appended different nouns related to a variety of occupations, such as “student”, “teacher”, and “lawyer”, to these words and used them as queries for Google image search. Subsequently, we collected images returned from the search engine and

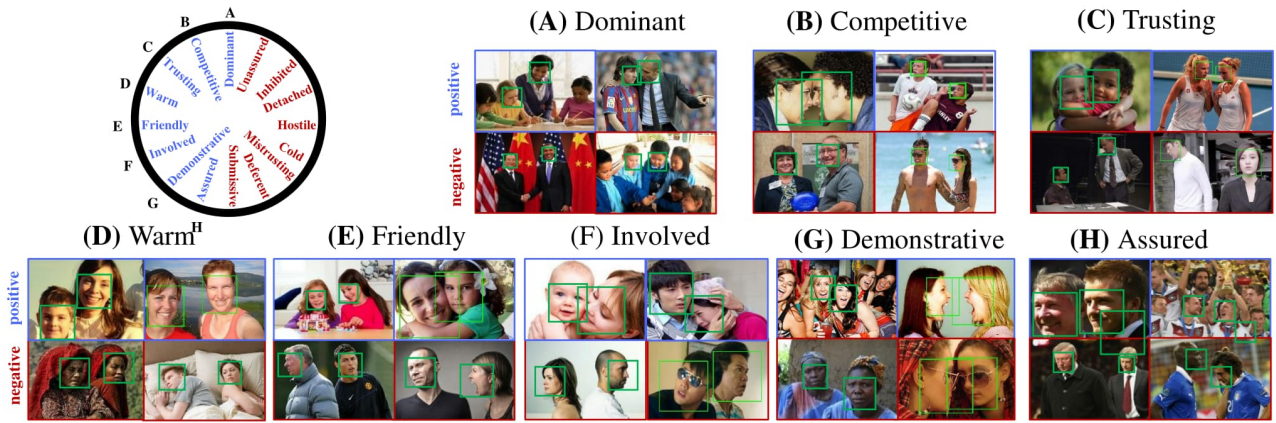
run a face detector [81] to obtain face regions from these images. Similar to other existing expression datasets [12, 20], each of the face images was manually annotated as one of the seven basic expression categories: “angry”, “disgust”, “fear”, “happy”, “sad”, “surprise”, or “neutral”. Non-face images were removed in the annotation process.

#### 3.2 Interpersonal Relation Dataset

To investigate the detectability of relation traits from a pair of face images, we built a new dataset containing 8,016 images chosen from web and movies. Each image was labeled with faces’ bounding boxes and their pairwise relations. This is the first face dataset annotated with interpersonal relation traits. It is challenging because of large face variations including poses, occlusions, and illuminations. In addition, the images exhibit rich relation traits from various sources including news photos of politicians, photos in social media, and video frames in movies, as shown in Fig. 3.

Before we collected for annotations, we first defined the interpersonal relation traits based on the interpersonal circle proposed by Kiesler [33] that commonly used in psychological studies, where human relations are divided into 16 segments as shown in Fig. 3. Each segment has its

**Fig. 2:** Example images of the proposed ExpW dataset.



**Fig. 3:** The 1982 Interpersonal Circle (upper left) is proposed by Donald J. Kiesler, and commonly used in psychological studies [33]. The 16 segments in the circle can be grouped into 8 relation traits. The traits are non-exclusive therefore can co-occur in an image. In this study, we investigate the detectability and quantification of these traits from computer vision point of view. (A)-(H) illustrate positive and negative examples of the eight relation traits.

**Table 2:** Descriptions of interpersonal relation traits based on the 1982 interpersonal circle [33].

Relation trait	Descriptions	Example pair
Dominant	one leads, directs, or controls the other / dominates the conversation / gives advices to the other	teacher & student
Competitive	hard and unsmiling / contest for advancement in power, fame, or wealth	people in a debate
Trusting	sincerely look at each other / no frowning or showing doubtful expression / not-on-guard about harm from each other	partners
Warm	speak in a gentle way / look relaxed / readily to show tender feelings	mother & baby
Friendly	work or act together / express sunny face / act in a polite way / be helpful	host & guest
Involved	engaged in physical interaction / involved with each other / not being alone or separated	lovers
Demonstrative	talk freely being unreserved in speech / readily to express the thoughts instead of keep silent / act emotionally	friends in a party
Assured	express to each other a feeling of bright and positive self-concept, instead of depressed or helpless	teammates

**Table 3:** Example adjectives for relation traits defined by Donald J. Kiesler [33].

Relation trait	Positive	Negative
dominant	controlling/leading/influencing/commanding/dictatorial	equal/matched/
competitive	critical/driven/enterprising	content/approving/flattering/respectful
trusting	unguarded/generous/innocent	mistrusting/suspicious/cunning/vigilant
warm	gentle/pardoning/soft/absolving	cold/strict/icy/harsh/cruel
friendly	cooperative/helpful/devoted	hostile/harmful/impolite/rude
involved	outgoing/attached/active/sociable	detached/distant/alof
demonstrative	talkative/casual/suggestive	mute/controlled/silent/unresponsive
assured	confident/cheerful/self-reliant/cocky	dependent/unassured/helpless/depressed

opposite side in the circle, such as “friendly and hostile”. Therefore, the 16 segments can be considered as eight binary relation traits, whose descriptions [33] and examples are given in Table 2. We also provide positive and negative visual samples for each relation in Fig. 3, showing that they are visually perceptible. For instance, “friendly” and “competitive” are easily separable because of the conflicting meanings. It is worth pointing out that some relations are close semantically, such as “friendly” and “trusting”. To accommodate such cases, we do not forcefully suppress any

one of these relations during prediction but allowing a pair of faces to have more than one relation.

Annotating relations is non-trivial and subjective by nature. We requested five performing arts students to label each relation for each face pair independently. A label was accepted if more than three annotations were consistent. The inconsistent samples were presented again to the five annotators to seek for consensus. To facilitate the annotation task, we also provided multiple cues to the annotators. First, to help them understand the definition of the relation traits,

**Table 4:** A summary of attributes annotated in AFLW [36], CelebA [48] and the proposed ExpW datasets, each of which contains 24,386, 202,599, and 91,793 face images, respectively.

Attributes	Gender	Pose					Expression										Age							
	gender	left profile	left	frontal	right	right profile	angry	disgust	fear	happy	sad	surprise	neutral	smiling	mouth opened	narrow eyes	young	goatee	no beard	sideburns	5 o'clock shadow	gray hair	bald	mustache
AFLW [36]	✓	✓	✓	✓	✓																			
CelebA [48]	✓													✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ExpW						✓	✓	✓	✓	✓	✓	✓												

we listed ten related adjectives (see Table 3 for examples) defined by [33] for the positive and negative samples on each relation trait, respectively. Multiple example images were also provided. Second, for image frames selected from movies, the annotators were asked to get familiar with the plot. The subtitles were presented during the labeling process. Third, we defined some measurable rules for the annotation of all relation traits. For example, if two people open their mouths, the relation trait of “demonstrative” is considered as positive; If a teacher is teaching his student, the “dominant” trait is considered as positive; A trait is defined as negative if the annotator cannot find any evidence to support its positive existence. The average Fleiss’ kappa of the eight relation traits annotation is 0.62, indicating substantial inter-rater agreement.

#### 4 Facial Expression and Attributes Recognition

The recognition of facial expression and other relevant attributes such as gender and age play a critical role in our relation prediction framework. In this study, we train a deep convolutional network end-to-end to map raw imagery pixels to a representation space and then perform expression and attribute prediction simultaneously. The joint learning of facial expression and attributes allows us to capture rich facial representation more effectively thus preparing a strong starting point for interpersonal relation learning.

##### 4.1 Problem Formulation and Overview

A natural way to learn a deep representation that captures multiple attributes is by training a multitask network that jointly predicts these attributes given a face image [92]. This can be implemented directly by introducing multiple supervisory tasks during the network training. In our problem training a multitask network, unfortunately, is non-trivial:

1. *Missing attribute labels* - As discussed in Section 1, face datasets that can cover all different kinds of attributes can hardly be found. The ExpW dataset collected by us, and the few popular face datasets such as AFLW [36]

and CelebA [48] contain subsets of attributes useful for our problem, but these subsets rarely overlap, as shown in Table 4. For instance, AFLW only contains gender and poses, while the ExpW dataset only has expressions. The many missing labels prevent us from ‘fully’ exploit an image since it is labeled with an attribute subset rather than a complete attribute set. The problem may also lead to sparsity in the supervisory signal and thus increase the convergence difficulty during training.

2. *Heterogeneous distribution* - These datasets were collected from different sources, therefore, exhibit vastly disparate statistical distributions. Specifically, the AFLW dataset contains face images gathered from Flickr that typically hosts high-quality photographs. Whereas the image quality in CelebA and ExpW is much lower and more diverse. Since these datasets are labeled with different sets of attributes, a direct joint training would bias each attribute to be trained by the corresponding labeled data alone, instead of benefiting from the existence of unlabeled images.

We propose a novel learning framework to mitigate the aforementioned problems. In general, given the training faces from multiple heterogeneous sources, we aim to train a deep convolutional network (DCN) that can predict the *union set* of attributes of these datasets (*i.e.* all attributes in Table 4). The training process is divided into two stages, as summarized in Algorithm 1. Further details of each stage are provided in Sec. 4.2 and Sec. 4.3.

1. *Network initialization* - Firstly, we initialize the parameters of our deep convolutional network by training it to minimize the classification error on the attributes despite the missing attribute labels in some samples.
2. *Alternating attribute propagation and face representation learning* - We fine-tune the network from the first stage via an alternating optimization process for obtaining a better face representation. The process is depicted in Fig. 4. In each iteration of the optimization, we extract the deep representation from each face, and compute the prior of attribute co-occurrence, based on which we perform attribute propagation to infer the missing attribute annotations as pseudo attribute labels in a MRF.

**Algorithm 1** Overview of the proposed framework.**Input:**

Multiple face image datasets with potentially non-overlapped attribute annotations.

**Output:**

Face representation that captures the union of the attributes from input datasets.

**Stage 1 Training:**

1: Initialize the network filters  $\mathbf{K}$  by maximizing Eqn. (1).

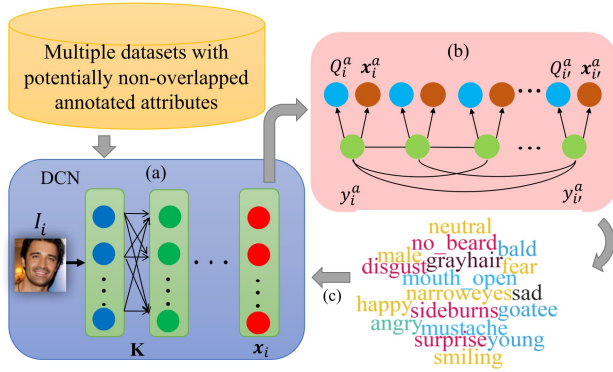
**Stage 2 Training:**

2: **for**  $m = 1$  to  $M$  **do**

3: Perform attribute propagation to fill up the missing labels by maximizing Eqn. (3).

4: Refine the network filters  $\mathbf{K}$  supervised by the ground truth and pseudo labels by minimizing the attribute classification error.

5: **end for**



**Fig. 4:** An illustration of the second stage training, in which we perform alternating optimization of representation learning and attribute propagation. (a) We extract face representation  $\mathbf{x}_i$  from the initialized DCN. (b) Given the face representation and attribute correlation, we perform attribute propagation in a Markov Random Field (MRF) to infer the missing attribute labels. (c) We refine the DCN by using the ground truth labels and pseudo labels generated from MRF.

We subsequently refine the network supervised by the ground truth attribute labels and newly generated pseudo attribute labels.

The second stage of Algorithm 1 helps to provide pseudo attribute labels that are missing initially for network fine-tuning. There are two advantages of this method: 1) The attribute propagation process does not require any prior knowledge of the problem at hand and thus can be applied given other datasets with an arbitrary number of missing labels. 2) Filling up the missing labels with pseudo labels naturally establish shared tasks among the datasets and gradually bridge the gap between datasets of different distributions. We show in the experiments (Sec. 6) that pseudo labels obtained in the attribute propagation step are crucial for good performance in the task of relation prediction.

## 4.2 First Stage: Network Initialization

The first stage of our training process is network initialization. Specifically, we first train the DCN (Fig. 4(a)) using a combined dataset, which includes AFLW, CelebA, and ExpW. Note that we do not perform attribute propagation at this stage but allow missing labels in samples.

Formally, let the network parameters be  $\mathbf{K}$ , an input face image  $\mathbf{I}_i$  is transformed to a higher level of representation represented as  $\mathbf{x}_i = \Phi(\mathbf{I}_i|\mathbf{K})$ , where  $\Phi(\mathbf{I}_i|\mathbf{K})$  denotes a nonlinear mapping parameterized by  $\mathbf{K}$ . We employ the Batch-Normalized Inception architecture presented in [29], where the network input is  $224 \times 224$  RGB image, and the generated face representation  $\mathbf{x}_i \in \mathbb{R}^{1024 \times 1}$ .

We assume the attributes are binary and thus we compute the probability for an attribute  $a$  by logistic regression. More precisely, given the attribute label  $y^a \in \{0, 1\}$ , we have  $p(y^a = 1|\mathbf{I}_i; \mathbf{K}, \mathbf{w}^a) = \frac{1}{1 + \exp(-\mathbf{w}^a \Phi(\mathbf{I}_i|\mathbf{K}))}$ , where  $\mathbf{w}^a$  are parameters of the logistic classifier. The network filter  $\mathbf{K}$  and classifier parameter  $\mathbf{w}^a$  can be obtained by maximizing the posterior probability:

$$\mathbf{K}^*, \mathbf{W}^{A*} = \underset{\mathbf{W}^A, \mathbf{K}}{\operatorname{argmax}} \sum_{i=1}^N \sum_{a=1}^{|A|} \log p(y_i^a | \mathbf{I}_i; \mathbf{w}^a, \mathbf{K}), \quad (1)$$

where  $N$  is the number of training samples,  $y_i^a$  is the ground truth label,  $A$  denotes the set of attributes, and  $\mathbf{W}^A = \{\mathbf{w}^a\}_{a \in A}$ . As a result, we can formulate a loss function with cross entropy for each attribute. The training process is conducted via back-propagation (BP) using stochastic gradient descent (SGD) [37].

Note that there are missing labels in the training set, which is combined from arbitrary datasets. To mitigate this issue, we mask the error of the missing attribute  $a$  of a training sample, and only back-propagate errors if the ground truth label of an attribute exists. Despite the missing labels, this simple approach provides a good initialization point for the second stage of the training.

## 4.3 Second Stage: Alternating Attribute Propagation and Face Representation Learning

**Formulation.** Following Algorithm 1, with the initialized network parameter  $\mathbf{K}$  and attribute classifier parameters  $\mathbf{W}^A$ , we subsequently perform attribute propagation to infer the missing attributes.

Attribute propagation is achieved based on two criteria: 1) Similarity of appearances between two faces, and 2) the correlation between attributes. The first criterion implies that the attributes of two faces are likely the same if their facial appearances are close to each other. The second criterion reflects the fact that some attributes, such as ‘happy’ and ‘smiling’, often co-occur.

With the above intuition, we formulate the attribute propagation problem in a MRF framework. In particular, as depicted in Fig. 4(b), each node in the MRF is an attribute label  $y_i^a$  for an image sample  $\mathbf{I}_i$ . Each edge describes the relation between the labels. For each node, we associate it with the observed variables  $\mathbf{x}_i$  representing the face representation obtained from the DCN, and  $Q_i^{a,a'}$ , which serves as a co-occurrence prior that indicates the tendency of an attribute  $a$  is present on a face  $i$ , given another attribute  $a'$  as condition.

We first provide the definition of the co-occurrence prior  $Q_i^{a,a'}$ . Given an attribute  $a$  and another attribute  $a' \in A \setminus a$ , we define  $Q_i^{a,a'}$  as

$$Q_i^{a,a'} = \begin{cases} \frac{\text{cov}(\mathbf{w}^a, \mathbf{w}^{a'})}{\sigma_{\mathbf{w}^a} \sigma_{\mathbf{w}^{a'}}} & \text{if } y_i^{a'} = 1 \text{ (positive)} \\ -\frac{\text{cov}(\mathbf{w}^a, \mathbf{w}^{a'})}{\sigma_{\mathbf{w}^a} \sigma_{\mathbf{w}^{a'}}} & \text{if } y_i^{a'} = 0 \text{ (negative)} \\ 0 & \text{if } y_i^{a'} \text{ is unlabeled.} \end{cases} \quad (2)$$

More precisely,  $Q_i^{a,a'}$  is assigned with the Pearson product-moment correlation coefficient [59], of which the sign is governed by the ground truth label of attribute  $a'$ , *i.e.*  $y_i^{a'}$ . The  $\text{cov}(\cdot)$  is the covariance, and  $\sigma$  is the standard deviation, while  $\mathbf{w}^a$  and  $\mathbf{w}^{a'}$  represent the parameters of the logistic classifier for the respective attribute. Intuitively, if attributes  $a$  and  $a'$  tend to co-occur, their  $\mathbf{w}^a$  and  $\mathbf{w}^{a'}$  are positively correlated. For instance, we have  $a = \text{“happy”}$ ,  $a' = \text{“smiling”}$ , and the Pearson correlation  $\frac{\text{cov}(\mathbf{w}^a, \mathbf{w}^{a'})}{\sigma_{\mathbf{w}^a} \sigma_{\mathbf{w}^{a'}}} = 0.3$ . For a face  $i$ , if the attribute “smiling” is annotated as positive (*i.e.*  $y_i^{a'} = 1$ ), then we have  $Q_i^{a,a'} = 0.3$ , suggesting that the “happy” attribute is present on the face given the “smiling” attribute. On the contrary, if the attribute “smiling” is absent (*i.e.*  $y_i^{a'} = 0$ ), then  $Q_i^{a,a'} = -0.3$ ,

---

**Algorithm 2** Alternating attribute propagation and face representation learning.

---

**Input:**

Face representation  $\mathbf{x}$ , and datasets with partially labeled attributes.

**Output:**

Pseudo label  $\mathbf{Y}^a$  on an attribute  $a$  for unlabeled data.

- 1: Compute the attribute co-occurrence prior  $\mathbf{Q}^a$  and extract face representation  $\mathbf{X}$ .
  - 2: For labeled data, use the original annotations; For unlabeled data, initialize the label by K-NN classification using the labeled data. Then we have the initial pseudo label  $\mathbf{Y}^a$ .
  - 3: Initialize the model parameter  $\Omega^a$  in Eqn. (4).
  - 4: Compute the affinity matrix  $V$  in Eqn. (6).
  - 5: Let iteration  $t = 0$ .
  - 6: **while** not converged **do**
  - 7:    $t = t + 1$ .
  - 8:   Infer a new  $\mathbf{Y}_t^a$  given the face representation  $\mathbf{X}$  and current model parameter  $\Omega_t^a$  (Eqn. (8)-(11)). Set  $\mathbf{Y}^a = \mathbf{Y}_t^a$ .
  - 9:   Update  $\Omega_t^a$  to maximize the log-likelihood of  $p(\mathbf{X}, \mathbf{Y}^a, \mathbf{Q}^a)$  by EM algorithm (Eqn. (13)-(14)).
  - 10: **end while**
- 

suggesting that the “happy” attribute is likely to absent too. We treat unannotated  $y_i^{a'}$  as a special case by forcing  $Q_i^{a,a'} = 0$ .

Let the face representation  $\mathbf{X} = \{\mathbf{x}_i\}$  and attribute co-occurrence prior  $\mathbf{Q}^a = \{Q_i^a\}$ , we maximize the following joint probability to obtain the attribute labels  $\mathbf{Y}^a = \{y_i^a\}$ :

$$\begin{aligned} p(\mathbf{X}, \mathbf{Y}^a, \mathbf{Q}^a) &= p(\mathbf{X}, \mathbf{Q}^a | \mathbf{Y}^a) p(\mathbf{Y}^a) \\ &= \frac{1}{Z} \prod_i \Phi(\mathbf{x}_i, Q_i^a | y_i^a) \prod_i \prod_{i' \in \mathcal{N}_i} \Psi(y_i^a, y_{i'}^a) \end{aligned} \quad (3)$$

where  $\Phi(\cdot)$ ,  $\Psi(\cdot)$  is the unary and pairwise term, respectively. The  $Z$  is the partition function, and  $\mathcal{N}_i$  denotes a set of face images, which are the neighbors of  $y_i^a$ .

We explain the unary and pairwise terms of Eqn. (3) as follows:

*Unary term* - We employ the Gaussian distribution to model the feature  $\mathbf{x}_i$  in the unary term  $\Phi(\cdot)$ . And we use the attribute co-occurrence prior as the prior probability. Specifically,

$$\Phi(\mathbf{x}_i, Q_i^a | y_i^a = \ell) \sim \mathcal{N}(\mathbf{x}_i | \mu_\ell^a, \Sigma_\ell^a) \cdot \prod_{a' \in A \setminus a} \mathcal{S}_{y_i^a}(Q_i^{a,a'}), \quad (4)$$

where  $\ell \in \{0, 1\}$ ,  $\mu_\ell$  and  $\Sigma_\ell$  denote the mean vector and covariance matrix of samples when  $y_i^a = \ell$ . Both  $\mu_\ell$  and  $\Sigma_\ell$  are obtained and updated during the inference process. For simplicity, we denote the model parameter  $\Omega^a = \{\mu_\ell^a, \Sigma_\ell^a\}$  in the following text. For  $\mathcal{S}_{y_i^a}(Q_i^{a,a'})$ , recall that given the attribute  $a'$ ,  $Q_i^{a,a'}$  denotes the prior that attribute  $a$  appears. Here we define  $\mathcal{S}_{y_i^a}$  as:

$$\mathcal{S}_{y_i^a}(Q_i^{a,a'}) = \begin{cases} \text{sigmoid}(Q_i^{a,a'}) & y_i^a = 1, \\ 1 - \text{sigmoid}(Q_i^{a,a'}) & y_i^a = 0. \end{cases} \quad (5)$$

Here “sigmoid” denotes a sigmoid transformation that maps the attribute co-occurrence prior from the range of  $[-1, 1]$  to  $[0, 1]$ . Hence,  $\mathcal{S}_{y_i^a}(Q_i^{a,a'})$  describes the prior that the attribute appears ( $y_i^a = 1$ ) or not ( $y_i^a = 0$ ).

*Pairwise term* - The pairwise term  $\Psi_p(\cdot)$  in Eqn.(3) is defined as

$$\Psi(y_i^a, y_{i'}^a) = \exp\{v_{ii'} \cdot \text{sign}(y_i^a, y_{i'}^a)\}, \quad (6)$$

where  $\text{sign}(\cdot)$  denotes a sign function:

$$\text{sign}(y_i^a, y_{i'}^a) = \begin{cases} 1 & y_i^a = y_{i'}^a \\ -1 & \text{otherwise.} \end{cases} \quad (7)$$

The variable  $v_{ii'}$  encodes the affinity between arbitrary pair of face image features  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$ . We obtain  $v_{ii'}$  via the spectral clustering approach presented in [90].



Firstly, we compute an affinity matrix  $V$  with entries  $v_{ii'} = \exp(-d^2(\mathbf{x}_i, \mathbf{x}_{i'})/\sigma_i\sigma_{i'})$  if  $\mathbf{x}_{i'}$  is within the  $h$ -nearest neighbors of  $\mathbf{x}_i$ , otherwise we set  $v_{ii'} = 0$ . We set  $h = 10$  in this study. The term  $d(\mathbf{x}_i, \mathbf{x}_{i'})$  is the  $\ell_2$ -distance between  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$ , and  $\sigma_i$  is the local scaling factor with  $\sigma_i = d(\mathbf{x}_i, \mathbf{x}_h)$ , where  $\mathbf{x}_h$  is the  $h$ -th nearest neighbor of  $\mathbf{x}_i$ . Then the normalized affinity matrix is obtained by  $V = D^{-\frac{1}{2}}VD^{-\frac{1}{2}}$ , where  $D$  is a diagonal matrix with  $D_{ii'} = \sum_{i'=1}^n v_{ii'}$ . Intuitively, Eqn. (6) penalizes face images with high affinity to be assigned with different attribute labels.

**Optimization.** Given the face representation  $\mathbf{X}$  and attribute co-occurrence prior  $\mathbf{Q}^a$ , we infer the missing attribute labels  $\mathbf{Y}^a$  by maximizing the joint probability of Eqn. (3).

Firstly, for the unlabeled data, we initialize the attribute  $a$  by K-NN classification in the space of  $\mathbf{x}$  using the labeled data. We keep the original attribute annotations for labeled data. Then we obtain  $\mu_\ell^a$  and  $\Sigma_\ell^a$  from the Gaussian of samples with  $y^a = \ell$ .

After the initialization of  $\Omega^a$ , we infer  $\mathbf{Y}^a$  and update the model parameter  $\Omega^a$  by repeating the following two steps in each optimization iteration  $t$ :

1. Infer a new  $\mathbf{Y}_t^a$  given the face representation  $\mathbf{X}$  and model parameter  $\Omega_t^a$ . Set  $\mathbf{Y}^a = \mathbf{Y}_t^a$ .
2. Given  $\mathbf{Y}^a$ , update  $\Omega_t^a$  to maximize the log-likelihood of  $p(\mathbf{X}, \mathbf{Y}^a, \mathbf{Q}^a)$  by Expectation-Maximization (EM) algorithm.

For the first step, we aim to obtain a new  $\mathbf{Y}_t^a$  given  $\mathbf{X}$  and model parameter  $\Omega_t^a$ . A natural way is to infer from the posterior:

$$p(\mathbf{Y}^a | \mathbf{X}, \mathbf{Q}^a, \Omega_t^a) = \frac{p(\mathbf{X}, \mathbf{Q}^a | \mathbf{Y}^a, \Omega_t^a)p(\mathbf{Y}^a)}{p(\mathbf{X}, \mathbf{Q}^a | \Omega_t^a)}. \quad (8)$$

However the computation of the term  $p(\mathbf{Y}^a)$  involves the interaction of each  $y_i^a$  and its neighborhood (*i.e.* the  $h$ -nearest neighbors in the space of  $\mathbf{x}$ ). Thus, it is intractable. Here we employ the mean field-like approximation [4] for  $p(\mathbf{Y}^a)$  computation, in which we assume each  $y_i^a$  is independent, and we set the value of its neighborhood  $\mathcal{N}_i$  constant when we compute  $p(y_i)$ . In this case, we have

$$p(\mathbf{Y}^a) = \prod_i p(y_i^a | \mathbb{Y}_{\mathcal{N}_i}^a), \quad (9)$$

where we denote the value of  $y_i$ 's neighborhood as  $\mathbb{Y}_{\mathcal{N}_i}^a \in \mathbb{R}^{|\mathcal{N}_i| \times 1}$ . For example, we can reuse the value in the previous iteration  $t-1$  (*i.e.*  $\mathbb{Y}_{\mathcal{N}_i}^a = \mathbf{Y}_{(t-1)\mathcal{N}_i}^a$ ). Because  $y_i^a \in \{0, 1\}$ , we have

$$\begin{aligned} p(y_i^a | \mathbb{Y}_{\mathcal{N}_i}^a) &= \frac{p(y_i, \mathbb{Y}_{\mathcal{N}_i}^a)}{\sum_{y_i \in \{0,1\}} p(y_i, \mathbb{Y}_{\mathcal{N}_i}^a)} \\ &= \frac{\frac{1}{Z} \prod_{j \in \mathcal{N}_i} \Psi(y_i^a, y_j^a)}{\sum_{y_i \in \{0,1\}} \frac{1}{Z} \prod_{j \in \mathcal{N}_i} \Psi(y_i^a, y_j^a)}. \end{aligned} \quad (10)$$

Since  $\mathbb{Y}_{\mathcal{N}_i}^a$  is fixed, the partition function  $Z$  is constant when we compute  $p(y_i, \mathbb{Y}_{\mathcal{N}_i}^a)$ . Thus  $Z$  can be eliminated in Eqn. (10). Combining Eqn. (4), Eqn. (8), Eqn. (10), we have

$$\begin{aligned} p(\mathbf{Y}^a | \mathbf{X}, \mathbf{Q}^a, \Omega_t^a) &= \prod_i p(y_i^a | \mathbb{Y}_{\mathcal{N}_i}^a, \mathbf{x}_i, Q_i^a, \Omega_t^a) \\ &= \prod_i \frac{\Phi(\mathbf{x}_i, Q_i^a | y_i^a, \Omega_t^a) \prod_{j \in \mathcal{N}_i} \Psi(y_i^a, y_j^a)}{\sum_{y_i^a \in \{0,1\}} \Phi(\mathbf{x}_i, Q_i^a | y_i^a, \Omega_t^a) \prod_{j \in \mathcal{N}_i} \Psi(y_i^a, y_j^a)}. \end{aligned} \quad (11)$$

Intuitively, the posterior  $p(y_i^a = \ell | \mathbb{Y}_{\mathcal{N}_i}^a, \mathbf{x}_i, \Omega_t^a)$  is proportional to the likelihood of setting  $y_i^a = \ell$ , with the neighborhood's value fixed. Then this posterior can be computed directly for each face  $i$ . To this end, we have

$$\mathbf{Y}_t^a = \{y_i^a\}_{i=1}^N \quad (12)$$

where for the unlabeled samples,  $y_i^a$  is simulated based on the posterior  $p(y_i^a = \ell | \mathbb{Y}_{\mathcal{N}_i}^a, \mathbf{x}_i, \Omega_t^a)$  (*i.e.* the probability of setting  $y_i^a = \ell$  is proportional to  $p(y_i^a = \ell | \mathbb{Y}_{\mathcal{N}_i}^a, \mathbf{x}_i, \Omega_t^a)$ ). For the annotated samples, we use the annotation directly.

For the second step, we aim to maximize the log-likelihood of  $p(\mathbf{X}, \mathbf{Y}^a, \mathbf{Q}^a)$  by updating the model parameter  $\Omega^a$  in an EM algorithm. Since  $\Omega^a = \{\mu^a, \Sigma^a\}$  only relates to  $\Phi(\mathbf{x}_i, Q_i^a | \Omega^a)$ , which is a Gaussian distribution, we update  $\Omega^a$  by

$$\mu_\ell^a = \frac{1}{|N_\ell|} \sum_{i \in N_\ell} \mathbf{x}_i, \quad (13)$$

$$\Sigma_\ell^a = \frac{1}{|N_\ell|} \sum_{i \in N_\ell} (\mathbf{x}_i - \mu_\ell^a) \cdot (\mathbf{x}_i - \mu_\ell^a)^\top, \quad (14)$$

where  $N_\ell$  denotes the subset of face images in which  $y_i^a = \ell$ .

The optimization of the above two steps ends when the posterior  $p(y_i^a | \mathbb{Y}_{\mathcal{N}_i}^a, \mathbf{x}_i, Q_i^a, \Omega_t^a)$  converged. The output attribute  $\mathbf{Y}^a$  is assigned with the final inferred  $\mathbf{Y}_t^a$ . Note that we use the original annotations. The optimization process is summarized in Algorithm 2.

## 5 Interpersonal Relation Prediction from Face Images

We have obtained a DCN that captures rich face representation through joint training with heterogeneous attribute sources. Next, we aim to jointly consider pairwise faces for interpersonal relation prediction.

We begin by arranging two identical DCNs obtained in Sec. 4 in a Siamese-like architecture as shown in Fig. 5. Using the interpersonal relation dataset introduced in Sec. 3.2, we train the new Siamese network end-to-end to map raw pixels of a pair of face images to relation traits.

As shown in Fig. 5, given an image with a detected pair of face, which is denoted as  $\mathbf{I}^r$  and  $\mathbf{I}^l$ , we extract high-level features  $\mathbf{x}^r$  and  $\mathbf{x}^l$  using two DCNs respectively. These two

DCNs have identical network structure as the one we use for expression recognition (see Sec. 4). Let  $\mathbf{K}^r$  and  $\mathbf{K}^l$  denote the network parameters. So we have  $\forall \mathbf{x}^r, \mathbf{x}^l \in \mathbb{R}^{1024 \times 1}$ . A weight matrix,  $\mathbf{W}^R \in \mathbb{R}^{2048 \times 256}$ , projects the concatenated feature vectors to a space of shared representation  $\mathbf{x}_g \in \mathbb{R}^{256}$ , which is utilized to predict a set of relation traits,  $\mathbf{g} = \{g_i\}_{i=1}^8, \forall g_i \in \{0, 1\}$ . Each relation is modeled as a single binary classification task, parameterized by a weight vector,  $\mathbf{w}_{g_i} \in \mathbb{R}^{256}$ .

In addition to the face images, we incorporate some spatial cues to train the deep network as shown in Fig. 5. The spatial cues include:

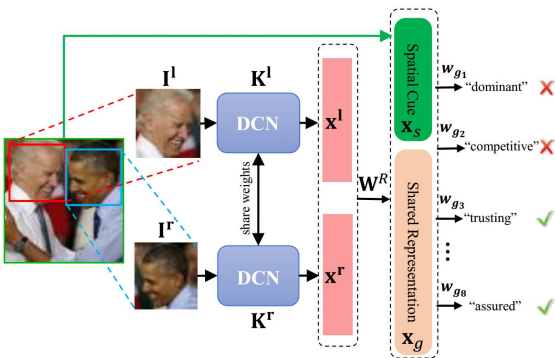
1. Two faces' positions  $\{x^l, y^l, w^l, h^l, x^r, y^r, w^r, h^r\}$ , representing the  $x$ - $y$ -coordinates of the upper-left corner, width, and height of the bounding boxes;  $w^l$  and  $w^r$  are normalized by the image width. Similar for  $h^l$  and  $h^r$
2. The relative faces' positions:  $\frac{x^l - x^r}{w^r}, \frac{y^l - y^r}{h^l}$
3. The ratio between the faces' scales:  $\frac{w^l}{w^r}$

The above spatial cues are concatenated as a vector,  $\mathbf{x}_s$ , and combined with the shared representation  $\mathbf{x}_g$  for learning relation traits.

Each binary variable  $g_i$  can be predicted by linear regression,

$$g_i = \mathbf{w}_{g_i}^T [\mathbf{x}_s; \mathbf{x}_g] + \epsilon, \quad (15)$$

where  $\epsilon$  is an additive error random variable, which is distributed following a standard logistic distribution,  $\epsilon \sim \text{Logistic}(0, 1)$ .  $[\cdot; \cdot]$  indicates the column-wise concatenation of two vectors. Therefore, the probability of  $g_i$  given  $\mathbf{x}_g$  and  $\mathbf{x}_s$  can be written as a sigmoid function,  $p(g_i = 1 | \mathbf{x}_g, \mathbf{x}_s) = 1 / (1 + \exp\{-\mathbf{w}_{g_i}^T [\mathbf{x}_s; \mathbf{x}_g]\})$ , indicating that  $p(g_i | \mathbf{x}_g, \mathbf{x}_s)$  is a Bernoulli distribution,  $p(g_i | \mathbf{x}_g, \mathbf{x}_s) = p(g_i = 1 | \mathbf{x}_g, \mathbf{x}_s)^{g_i} (1 - p(g_i = 1 | \mathbf{x}_g, \mathbf{x}_s))^{1-g_i}$ .



**Fig. 5:** Overview of the network for interpersonal relation learning. The input is two face images and we extract the representation by two identical DCN, which is initialized by learning on multiple attribute datasets (see Sec. 4). Then we perform relation traits reasoning using face representation and additional spatial cues. The output is eight binary values that encode the different dimensions of relation traits.

In addition, the probabilities of  $\mathbf{w}_{g_i}$ ,  $\mathbf{W}^R$ ,  $\mathbf{K}^l$ , and  $\mathbf{K}^r$  can be modeled by the standard normal distributions. For example, suppose  $\mathbf{K}$  contains  $K$  filters, then  $p(\mathbf{K}) = \prod_{j=1}^K p(\mathbf{k}_j) = \prod_{j=1}^K \mathcal{N}(\mathbf{0}, \mathcal{I})$ , where  $\mathbf{0}$  and  $\mathcal{I}$  are an all-zero vector and an identity matrix respectively, implying that the  $K$  filters are independent. Similarly, we have  $p(\mathbf{w}_{g_i}) = \mathcal{N}(\mathbf{0}, \mathcal{I})$ . Furthermore,  $\mathbf{W}^R$  can be initialized by a standard matrix normal distribution [22], *i.e.*,  $p(\mathbf{W}^R) \propto \exp\{-\frac{1}{2} \text{tr}(\mathbf{W}^R \mathbf{W}^{R^T})\}$ , where  $\text{tr}(\cdot)$  indicates the trace of a matrix.

Combining the above probabilistic definitions, the deep network is trained by maximizing a posterior probability,

$$\begin{aligned} \arg \max_{\Omega} p(\{\mathbf{w}_{g_i}\}_{i=1}^8, \mathbf{W}, \mathbf{K}^l, \mathbf{K}^r | \mathbf{g}, \mathbf{x}_g, \mathbf{x}_s, \mathbf{I}^l, \mathbf{I}^r) \propto \\ \left( \prod_{i=1}^8 p(g_i | \mathbf{x}_g, \mathbf{x}_s) p(\mathbf{w}_{g_i}) \right) \left( \prod_{j=1}^K p(\mathbf{k}_j^l) p(\mathbf{k}_j^r) \right) p(\mathbf{W}^R), \\ \text{s.t. } \mathbf{K}^r = \mathbf{K}^l \end{aligned} \quad (16)$$

where  $\Omega = \{\{\mathbf{w}_{g_i}\}_{i=1}^8, \mathbf{W}^R, \mathbf{K}^l, \mathbf{K}^r\}$  and the constraint means the filters are tied. Note that  $\mathbf{x}_g$  and  $\mathbf{x}_s$  represent the hidden features and the spatial cues extracted from the left and right face images, respectively. Thus, the variable  $g_i$  is independent with  $\mathbf{I}^l$  and  $\mathbf{I}^r$ , given  $\mathbf{x}_g$  and  $\mathbf{x}_s$ .

By taking the negative logarithm of Eqn.(16), it is equivalent to minimizing the following loss function

$$\begin{aligned} \arg \min_{\Omega} \sum_{i=1}^8 \left\{ \mathbf{w}_{g_i}^T \mathbf{w}_{g_i} - (1 - g_i) \ln(1 - p(g_i = 1 | \mathbf{x}_g, \mathbf{x}_s)) - \right. \\ \left. g_i \ln p(g_i = 1 | \mathbf{x}_g, \mathbf{x}_s) \right\} + \sum_{j=1}^K (\mathbf{k}_j^r{}^T \mathbf{k}_j^r + \mathbf{k}_j^l{}^T \mathbf{k}_j^l) + \text{tr}(\mathbf{W}^R \mathbf{W}^{R^T}), \\ \text{s.t. } \mathbf{k}_j^r = \mathbf{k}_j^l, j = 1 \dots K \end{aligned} \quad (17)$$

where the second and the third terms correspond to the traditional cross-entropy loss, while the remaining terms indicate the weight decays [53] of the parameters. Equation (17) is defined over single training sample and is a highly nonlinear function because of the hidden features  $\mathbf{x}_g$ . Here we first initialize  $\mathbf{K}^l$  and  $\mathbf{K}^r$  by the representation we learn in Sec. 4. Then Eqn. (17) is solved by stochastic gradient descent [37].

## 6 Experiments

We divide our experiments into two subsections. Section 6.1 examines the effectiveness of our base DCN on facial expression and attributes recognition. Section 6.2 evaluates our full Siamese framework for interpersonal relation prediction.

## 6.1 Facial Expression and Attributes Recognition

**Dataset.** We evaluated our base DCN on the combined dataset of AFLW, CelebA, and ExpW. From the total of 318,778 face images, we selected 5,400 images for testing and the remaining were reserved for training and validation. The test images consisted of 3,000 CelebA, 1,000 AFLW, and 1,400 ExpW images. We ensured that the ExpW test partition was balanced in their seven facial expression classes, *i.e.* all expression class had 200 samples. Note that this rule was not enforced in other attribute categories.

In addition to this combined dataset, we also evaluated our approach on the Static Facial Expressions in the Wild (SFEW) dataset [12] and CK+ [49] datasets.

**Evaluation metric.** To account for the imbalanced positive and negative attribute samples, a balanced accuracy is adopted as the evaluation metric:

$$accuracy = 0.5 \times (n_p/N_p + n_n/N_n), \quad (18)$$

where  $N_p$  and  $N_n$  are the numbers of positive and negative samples, while  $n_p$  and  $n_n$  are the numbers of true positive and true negative.

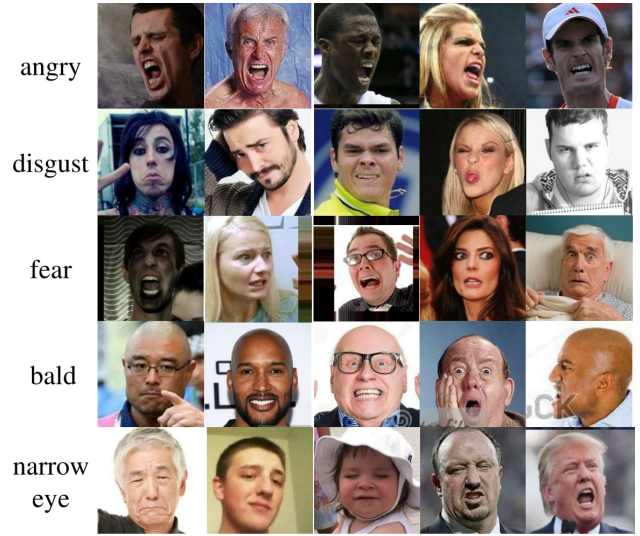
**Implementation.** We implemented the proposed deep model with MXNet [72] library. Data augmentation by random translation and mirroring were introduced in the training process. The mini-batch size was fixed to 32, and the learning rate was 0.001 with a momentum rate of 0.9. Following Algorithm 1, the first initialization stage took 30 epochs to converge, while the second stage on attribute propagation consumed another 10 epochs (*i.e.*,  $M = 10$ ).

**Results on the combined AFLW, CelebA, and ExpW.** We trained two variants of our DCN using the combined dataset:

1. Baseline DCN - it is trained without both attribute propagation.
2. DCN+AP - it is trained with attribute propagation (*i.e.* full model).

For completeness, we additionally trained a baseline classifier by extracting HOG features from the given face images, and we used a linear support vector machine (SVM) to train a binary classifier (*i.e.*, HOG+SVM) for each attribute. In the SVM learning process, we adjusted the weight of each class as inversely proportional to the class frequency in the training data. This helped in mitigating the imbalanced class issue.

The balanced accuracy of each method is reported in Table 5. It is observed that in general, attribute propagation helps, especially on attributes with rare positive samples such as “narrow eyes” and “goatee”. We conjecture that attribute propagation allows the proposed model to effectively leverage samples from multiple datasets, which are not annotated initially.



**Fig. 6:** Examples of automatically annotated positive attribute examples via the proposed attribute propagation (discussed in Sec. 4.3).

To further compare with existing attribute recognition methods, we follow the training and testing splits of CelebA [48] (as for AFLW and ExpW, we use the same training data as the previous experiments). The performance is summarized in Table 6. Note that we follow the convention of [48], and use the overall classification accuracy instead of the balanced accuracy as Eqn. (18). We can observe that by fusing multiple datasets, our proposed method achieves superior performance compared to state-of-the-art methods.

In Table 7, we show the average balanced accuracy over different iterations of the alternating attribute propagation and representation learning process (see Sec. 4.3). The gradually improved accuracy over iterations demonstrates that the alternating optimization process is beneficial. Figure 6 shows a few initially unlabeled positive attribute samples that are automatically annotated via attribute propagation. It is worth pointing out that many of this unlabeled samples are challenging in terms of their unconstrained poses and expressions.

**Expression Recognition on SFEW [12].** To demonstrate the effectiveness of the proposed DCN for facial expression recognition, we evaluated its performance on the challenging Static Facial Expressions in the Wild (SFEW) 2.0 dataset [12]. The dataset is a static subset of Acted Facial Expressions in the Wild (AFEW) dataset [12], which captures natural and versatile expressions from movies. Since the label for the test set is not publicly available, we follow the training/validation splits of the released dataset, we evaluated two variants of our method: 1) Our trained DCN+AP without fine-tuning on SFEW training partition, and 2) Our trained full model DCN+AP with fine-tuning on SFEW

**Table 5:** Balanced accuracies (%) over different attributes.

Attributes	Gender		Pose					Expression								Age									
	average	gender	left profile	left	frontal	right	right profile	angry	disgust	fear	happy	sad	surprise	neutral	smiling	mouth opened	narrow eyes	young	goatee	no beard	sideburns	5 o'clock shadow	gray hair	bald	mustache
HOG+SVM	71.0	83.2	73.8	65.7	88.3	60.3	70.1	54.3	54.8	56.2	71.3	58.4	61.2	68.4	84.5	79.7	56.3	72.9	75.6	88.4	75.8	72.4	85.9	75.4	70.4
Baseline DCN	76.1	96.5	75.0	56.3	87.3	51.8	74.2	63.5	50.0	50.0	81.9	64.0	71.0	75.0	93.0	94.2	63.3	<b>84.4</b>	84.8	92.8	88.6	82.8	87.7	86.1	73.4
DCN+AP	<b>80.9</b>	<b>97.0</b>	<b>78.4</b>	<b>67.3</b>	<b>90.0</b>	<b>62.1</b>	<b>77.9</b>	<b>72.1</b>	<b>56.5</b>	<b>58.7</b>	<b>83.8</b>	<b>69.1</b>	<b>74.2</b>	<b>76.0</b>	<b>93.3</b>	<b>94.5</b>	<b>73.5</b>	83.5	<b>90.1</b>	<b>92.5</b>	<b>92.5</b>	<b>88.3</b>	<b>92.2</b>	<b>93.0</b>	<b>83.9</b>

**Table 6:** Attribute Recognition Accuracy on CelebA.

Method	smiling	mouth opened	narrow eyes	young	goatee	no beard	sideburns	5 o'clock shadow	gray hair	bald	mustache	gender
FaceTracer [38]	89	87	82	80	93	90	94	85	90	89	91	91
PANDA-w [91]	89	82	79	77	86	87	90	82	88	92	83	93
PANDA-l [91]	92	93	84	84	93	93	93	88	94	96	93	97
Liu <i>et al.</i> [48]	92	92	81	87	95	95	96	91	97	98	95	<b>98</b>
MCNN-AUX [15]	93	94	87	88	97	96	98	95	98	<b>99</b>	<b>97</b>	<b>98</b>
ours	<b>94</b>	<b>95</b>	<b>89</b>	<b>91</b>	<b>98</b>	<b>97</b>	<b>98</b>	<b>96</b>	<b>99</b>	<b>99</b>	<b>98</b>	<b>98</b>

**Table 7:** Average balanced accuracies (%) over different iterations of the alternating attribute propagation and representation learning process.

Iteration	M=1	M=3	M=5	M=7	M=9	M=10
Accuracy	78.4	79.2	79.3	79.8	80.8	80.9

training partition. Our model treats each expression as a binary attribute, the expression with the highest predicted probability is selected as the classification result.

We compared our method with the following approaches:

1. PHOG+LPQ [12] - the Pyramid of Histogram of Gradients (PHOG) and Local Phase Quantization (LPQ) [11] are computed and concatenated to form the feature of a face, and a non-linear SVM is used for expression classification.
2. MBP [41] - expression recognition with Mapped Binary Patterns (MBP), which is proposed in [41].
3. AU-Aware Features [86] - expression recognition by exploiting facial action-unit aware features.
4. Microsoft Emotion API [1] - emotion API of Microsoft cognitive services. Since it is a commercial API, we use the service directly without fine-tuning in on the SFEW training partition.
5. DCN of [54] - AlexNet [37] pretrained with ImageNet [65] and FER [20] datasets, and finetune on the SFEW training dataset.
6. DCN of [88] - A customized DCN (five convolutional layers and two fully connected layers) pretrained with

**Table 8:** Accuracies on the validation set of SFEW dataset [12].

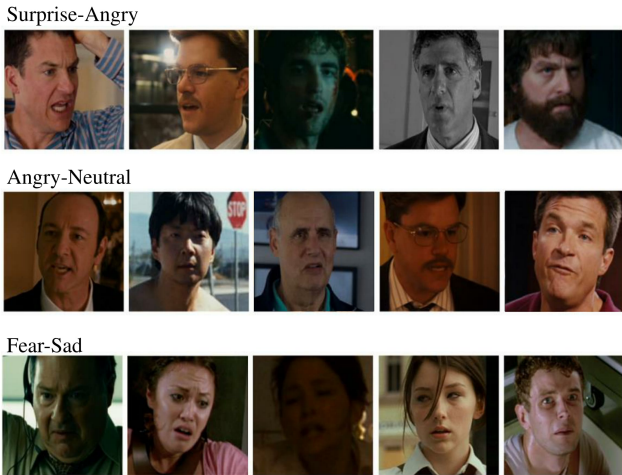
Method	Training/Fine-tuning on SFEW	Accuracy
PHOG+LPQ [12]	yes	35.93%
MBP [41]	yes	41.92%
AU-Aware features [86]	yes	44.04%
Microsoft Emotion API [1]	no	47.71%
DCN of [54]	yes	48.50%
Single DCN of [88]	yes	52.29%
Ensemble DCNs of [88] <sup>1</sup>	yes	55.96%
Our Baseline DCN	no	45.51%
Our DCN+AP	no	49.77%
Our Baseline DCN	yes	52.06%
Our DCN+AP	yes	55.27%

<sup>1</sup> This result is obtained from an ensemble of five DCNs.

FER [20] dataset, and fine-tune on the SFEW training dataset.

Table 8 summarizes the performances of various approaches evaluated on the SFEW dataset. Following the convention of current studies [88,43,32,46], we use the overall classification accuracy instead of the balanced accuracy as Eqn. (18). Our approach, with and without fine-tuning on SFEW training partition, outperforms state-of-the-art methods. Again, it is observed that our model is benefited from alternating optimization with attribute propagation. Figure 7 shows some failure cases. Most errors were caused by ambiguous cases.

**Expression Recognition on CK+ [49].** For completeness, we also evaluated our method on CK+ [49] since it is a classic dataset for expression recognition. CK+ contains 327 image sequences where each sequence presents a face with gradual expression evolution from a neutral to a peak facial expression. Each sequence is annotated with one of the six prototypical expressions, *i.e.*, angry, happy, surprise, sad, disgust, fear, or a non-standard expression (*i.e.* contempt). Following the widely used evaluation protocol [46,32,97], we selected the last three frames of each sequence for training/testing purpose. The first frame of each sequence was regarded as the “neutral” expression. Consequently, we obtained 1,308 images for our 10-fold cross-validation. The face identity in each fold was remained exclusive. As in the



**Fig. 7:** Example of failure cases of our approach (DCN+AP) on the SFEW validation set. The text above each row denotes the ground truth and predicted result, e.g., “Surprise-Angry” means the surprise expression is misclassified as angry. Most failures were caused by ambiguity in facial expressions.

**Table 9:** Accuracies on the CK+ dataset [49] with six prototypical facial expressions.

Method	Accuracy
CSPL [98]	89.9%
LBPSVM[67]	95.1%
BDBN [46]	96.7%
PPDN [97]	97.3%
Zero-bias CNN [32]	98.3%
Our Method	98.9%

SFEW experiments, we fine-tuned our trained DCN+AP on the training samples of each fold.

Table 9 presents the comparative results of our method and other state-of-the-arts. To be consistent with other methods, the averaged accuracy of the six basic expressions are reported. Similar to our approach, BDBN [46], PPDN [97], and Zero-bias CNN [32] also adopted different kinds of deep networks. Our approach still achieves better result although the performance on CK+ is nearly saturated.

## 6.2 Interpersonal Relation Prediction

**Dataset.** The evaluation of interpersonal relation learning was performed on the dataset described in Sec. 3.2. We divided the dataset into training and test partitions of 7,226 and 790 images, respectively. The face pairs in these two partitions were mutually exclusive, containing no overlapped identities. Table 10 presents the statistics of this dataset.

**Evaluation metric.** We adopt the same balanced accuracy in Eqn. (18).

**Table 10:** Statistics of the interpersonal relation dataset.

Relation trait	training		testing	
	#positive	#negative	#positive	#negative
dominant	418	6808	112	678
competitive	344	6882	70	720
trusting	6261	965	606	184
warm	6176	1050	615	175
friendly	6733	493	728	62
involved	6360	866	686	104
demonstrative	6494	732	689	101
assured	6538	688	673	117

**Baselines.** As discussed in Sec. 5, our full model combines the two DCNs pre-trained for expression and attribute recognition in a Siamese-like architecture, as shown in Fig. 5. We call this model as “S-DCN”.

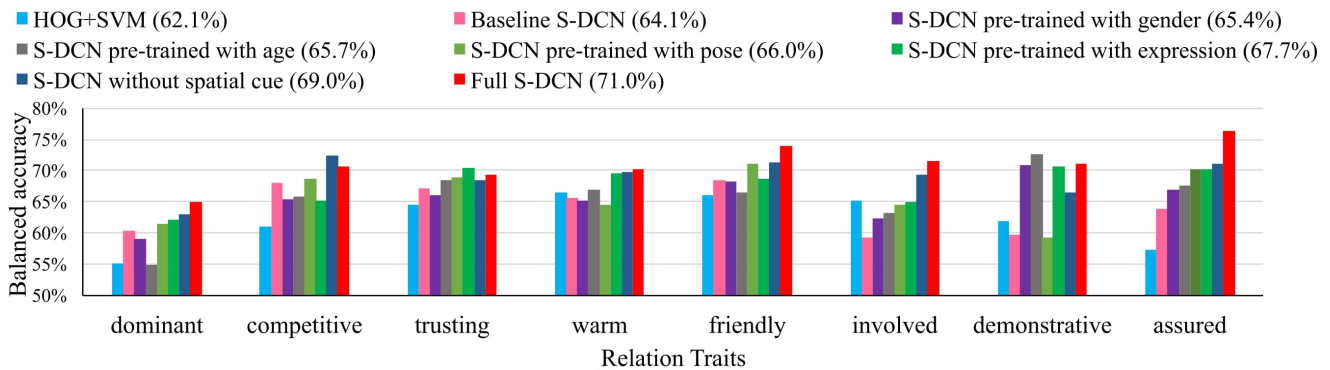
We evaluated several variants of this network.

1. Baseline S-DCN - We trained a model similar to S-DCN in Fig. 5, but without using the DCN pre-trained for expression and attribute recognition. Instead, the parameters of the two DCNs were randomly initialized.
2. S-DCN with its DCN pre-trained with selected attributes - To examine the influences of different attribute groups, we pre-trained four DCN variants using only one group of attribute (*i.e.*, expression, age, gender, and pose), respectively.
3. S-DCN without spatial cue - We trained a S-DCN with DCN pre-trained with all the attributes but the spatial cue (discussed in Sec. 5) was not used.
4. Full S-DCN - We trained a S-DCN with DCN pre-trained with all the attributes and used the spatial cue as discussed in Sec. 5.

In addition, we established a baseline “HOG+SVM” - we extracted the HOG features from the given face images. The features from two faces were then concatenated and a linear support vector machine (SVM) was employed to train a binary classifier for each relation trait.

**Results.** Figure 8 shows the accuracies of different variants. All variants of the proposed S-DCN outperform the baseline HOG+SVM. We observe that the cross-dataset expression and attribute pre-training is beneficial since pre-training with any of the attribute groups improves the overall performance. In particular, pre-training with expression attributes outperforms other groups of attributes (improving from 64.1% to 67.7%). This is not surprising since interpersonal relation is largely manifested from expression. The pose attributes come next in terms of influence to relation prediction. The result is also expected since when people are in a close or friendly relation, they tend to look at the same direction or face each other.

Finally, the spatial cue is shown to be useful for relation prediction. However, we also observe that not every trait



**Fig. 8:** Relation traits prediction performance. The number in the legend indicates the average accuracy of the according method across all the relation traits.

**Table 11:** Balanced accuracies (%) on the movie testing subset.

Method	Balanced Accuracy
HOG+SVM	59.22%
Baseline S-DCN	62.42%
S-DCN (DCN pre-trained with gender)	63.10%
S-DCN (DCN pre-trained with age)	64.67%
S-DCN (DCN pre-trained with pose)	62.83%
S-DCN (DCN pre-trained with expression)	65.36%
S-DCN without spatial cue	68.17%
Full S-DCN	70.20%

is improved by the spatial cue and some are degraded. This is because currently we simply use the face scale and location directly, of which the distribution is inconsistent in images from different sources. For example, some close-shot photographs may be used to show competing people and their expression in detail, while in some movies, competing people may stand far away from each other. As for the relation traits, “dominant” is the most difficult trait to predict as it needs to be determined by more complicated factors, such as one’s social role and the environmental context.

To factor out any potential subjective judgement arisen from the data annotation process, we evaluated S-DCN on a subset of 522 movie frames extracted from the test data. This subset is more ‘objective’ since annotators were provided with richer auxiliary cues for relation annotation. Table 11 shows the average balanced accuracy on the eight relation traits of the baseline and the variants of the proposed S-DCN. The results further suggest the reliability of the proposed approach.

Some positive and negative predictions on different relation traits are shown in Fig. 9(a). It can be observed that the proposed approach is capable of handling images in different scenes and faces with large expression variations. We show some false positives in Fig. 9(b), which are partly caused by the lack of context. For example, in the first image of Fig. 9(b), the two characters were having a serious

conversation. The algorithm had no access to the context that they were reading a book and thus guessed that they were competing. Our method also failed given faces with a large degree of occlusions.

More qualitative results are presented in Fig. 10. Positive relation traits, such as “trusting”, “warm”, “friendly” are inferred between the US President *Barack Obama* and his family members. Interestingly, “dominant” trait is predicted between him and his daughter (Fig. 10 (b)). Fig. 10(c) includes the image for *Angela Merkel*, Chancellor of Germany, which is usually used in the news articles on US spying scandal, showing a low tendency on the “trusting” trait, while a high tendency on the “competitive” trait. This relation is quite different from that of Fig. 10(d), where *Obama* and the British Prime Minister *David Cameron* were watching a basketball game.

We show an example of application of using our method to automatically profile the relations among the characters in a movie. We chose the movie *Iron Man* and focused on different interaction patterns, such as conversation and conflict, of the main roles “*Tony Stark*” and “*Pepper Potts*”. Firstly, we applied a face detector to the movie and selected those frames that captured the two roles. Then, we applied our approach on each frame to infer their relation traits. The predicted probabilities were averaged across 5 neighboring frames to obtain a smooth profile. Figure 11 shows a video segment with the traits of “friendly” and “competitive”. Our method accurately captures the friendly talking scene and the moment when Tony and Pepper were in a conflict, where the “competitive” trait is assigned with a high probability while the “friendly” trait is low.

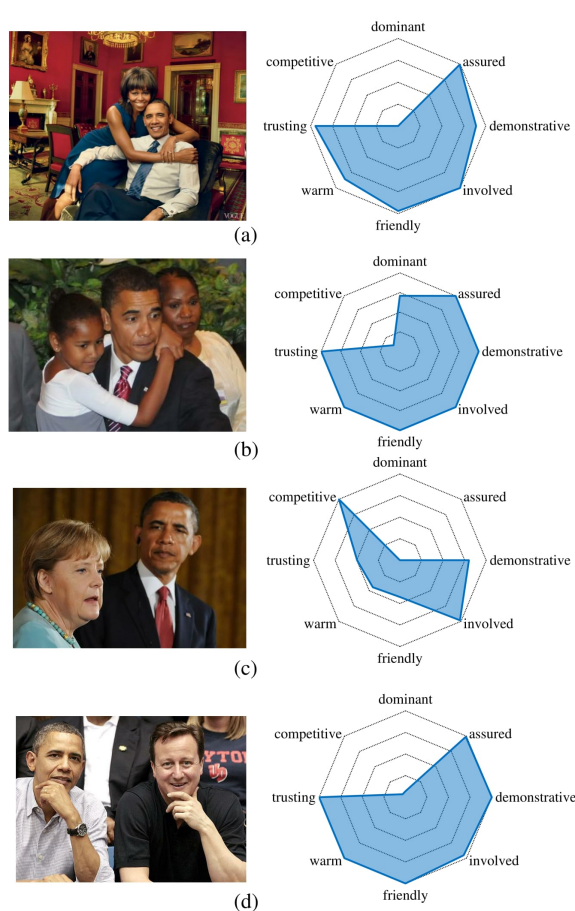
## 7 Conclusion

In this work, we studied a new challenging problem of predicting interpersonal relation from face images. We decomposed our solution into two steps. We began with



**Fig. 9:** (a) Correct positive and negative prediction results on different relation traits. (b) False positives on “competitive”, “assured” and “demonstrative” relation traits (from left to right).

training a reliable deep convolutional network for recognizing facial expression and rich attributes (gender, age, and pose) from single face images. We addressed the problem of learning from heterogeneous data sources with potentially missing attribute labels. This was achieved through a novel approach that leverages the inherent correspondences among heterogeneous sources by attribute propagation in a graphical model. Initialized by the deep convolutional network learned in the first step, a Siamese-like framework



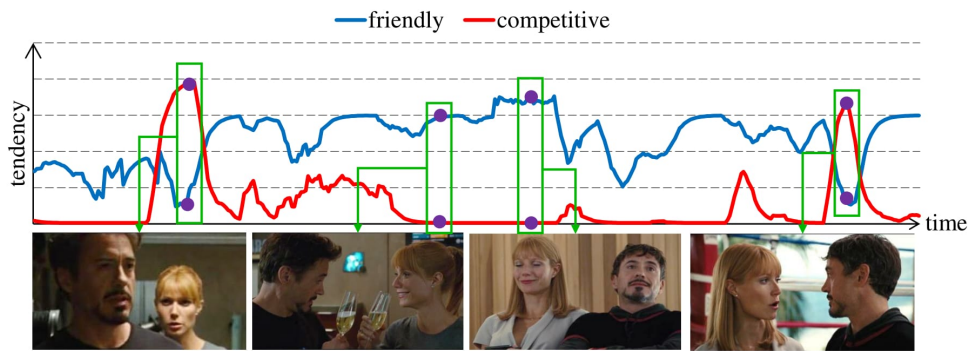
**Fig. 10:** The relation traits predicted by our full model with spatial cue (Full S-DCN). The polar graph beside each image indicates the tendency for each trait to be positive.

is proposed to learn an end-to-end mapping from raw pixels of a pair of face images to relation traits. Extensive experiments demonstrate the effectiveness of the proposed methods on facial expression recognition and interpersonal relation prediction. Future work will combine the face-based relation traits with body-driven immediacy cues [7] for more accurate interpersonal relation prediction.

**Acknowledgements** This work is supported by SenseTime Group Limited and the General Research Fund sponsored by the Research Grants Council of the Hong Kong SAR (CUHK 14241716, 14224316, 14209217).

**References**

1. Microsoft cognitive services (2016). URL <https://www.microsoft.com/cognitive-services/en-us/emotion-api>
2. Bi, W., Kwok, J.T.: Multilabel classification with label correlations and missing labels. In: AAAI Conference on Artificial Intelligence, pp. 1680–1686 (2014)



**Fig. 11:** Prediction for relation traits of “friendly” and “competitive” for the movie *Iron Man*. The probability indicates the tendency for the trait to be positive. It shows that the proposed approach can capture the friendly talking scene and the moment of conflict.

3. Bromley, J., Guyon, I., Lecun, Y., Säckinger, E., Shah, R.: Signature verification using a siamese time delay neural network. In: *Advances in Neural Information Processing Systems* (1994)
4. Celeux, G., Forbes, F., Peyrard, N.: EM procedures using mean field-like approximations for markov model-based image segmentation. *Pattern recognition* **36**(1), 131–144 (2003)
5. Chakraborty, I., Cheng, H., Javed, O.: 3D visual proxemics: Recognizing human interactions in 3D from a single image. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3406–3413 (2013)
6. Chen, Y.Y., Hsu, W.H., Liao, H.Y.M.: Discovering informative social subgraphs and predicting pairwise relationships from group photos. In: *ACM Multimedia*, pp. 669–678 (2012)
7. Chu, X., Ouyang, W., Yang, W., Wang, X.: Multi-task recurrent neural network for immediacy prediction. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3352–3360 (2015)
8. Cristani, M., Raghavendra, R., Del Bue, A., Murino, V.: Human behavior analysis in video surveillance: A social signal processing perspective. *Neurocomputing* **100**, 86–97 (2013)
9. Dahmane, M., Meunier, J.: Emotion recognition using dynamic grid-based hog features. In: *IEEE International Conference on Automatic Face & Gesture Recognition*, pp. 884–888 (2011)
10. Deng, Z., Vahdat, A., Hu, H., Mori, G.: Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
11. Dhall, A., Asthana, A., Goecke, R., Gedeon, T.: Emotion recognition using phog and lpq features. In: *IEEE International Conference on Automatic Face & Gesture Recognition and Workshops*, pp. 878–883 (2011)
12. Dhall, A., Ramana Murthy, O., Goecke, R., Joshi, J., Gedeon, T.: Video and image based emotion recognition challenges in the wild: EmotiW 2015. In: *ACM International Conference on Multimodal Interaction*, pp. 423–426 (2015)
13. Ding, L., Yilmaz, A.: Learning relations among movie characters: A social network perspective. In: *European Conference on Computer Vision* (2010)
14. Ding, L., Yilmaz, A.: Inferring social relations from visual concepts. In: *IEEE International Conference on Computer Vision*, pp. 699–706 (2011)
15. Emily M. Hand, R.C.: Attributes for improved attributes: a multi-task network utilizing implicit and explicit relationships for facial attribute classification. In: *AAAI Conference on Artificial Intelligence* (2017)
16. Fabian Benitez-Quiroz, C., Srinivasan, R., Martinez, A.M.: Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
17. Fathi, A., Hodgins, J.K., Rehg, J.M.: Social interactions: A first-person perspective. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2012)
18. Gallagher, A.C., Chen, T.: Understanding images of groups of people. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 256–263. *IEEE* (2009)
19. Girard, J.M.: Perceptions of interpersonal behavior are influenced by gender, facial expression intensity, and head pose. In: *ACM International Conference on Multimodal Interaction*, pp. 394–398 (2014)
20. Goodfellow, I., Erhan, D., Carrier, P.L., Courville, A., Mirza, et al.: Challenges in representation learning: A report on three machine learning contests (2013). URL <http://arxiv.org/abs/1307.0414>
21. Gottman, J., Levenson, R., Woodin, E.: Facial expressions during marital conflict. *Journal of Family Communication* **1**(1), 37–57 (2001)
22. Gupta, A.K., Nagar, D.K.: *Matrix variate distributions*. CRC Press (1999)
23. Hess, U., Blairy, S., Kleck, R.E.: The influence of facial emotion displays, gender, and ethnicity on judgments of dominance and affiliation. *Journal of Nonverbal Behavior* **24**(4), 265–283 (2000)
24. Hoai, M., Zisserman, A.: Talking heads: detecting humans and recognizing their interactions. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2014)
25. Hu, Y., Zeng, Z., Yin, L., Wei, X., Zhou, X., Huang, T.S.: Multi-view facial expression recognition. In: *IEEE International Conference on Automatic Face & Gesture Recognition*, pp. 1–6 (2008). DOI 10.1109/AFGR.2008.4813445
26. Huang, C., Li, Y., Loy, C.C., Tang, X.: Learning deep representation for imbalanced classification. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
27. Hung, H., Jayagopi, D., Yeo, C., Friedland, G., Ba, S., Odobez, J.M., Ramchandran, K., Mirghafori, N., Gatica-Perez, D.: Using audio and video features to classify the most dominant person in a group meeting. In: *ACM Multimedia* (2007)
28. Ibrahim, M., Muralidharan, S., Deng, Z., Vahdat, A., Mori, G.: A hierarchical deep temporal model for group activity recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
29. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: *Proceedings of The 32nd International Conference on Machine Learning*, pp. 448–456 (2015)



30. Joo, J., Li, W., Steen, F., Zhu, S.C.: Visual persuasion: Inferring communicative intents of images. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 216–223 (2014)
31. Jung, H., Lee, S., Yim, J., Park, S., Kim, J.: Joint fine-tuning in deep neural networks for facial expression recognition. In: *IEEE International Conference on Computer Vision* (2015)
32. Khorrami, P., Paine, T., Huang, T.: Do deep neural networks learn facial action units when doing expression recognition? In: *IEEE International Conference on Computer Vision Workshop* (2015)
33. Kiesler, D.J.: The 1982 interpersonal circle: A taxonomy for complementarity in human transactions. *Psychological Review* **90**(3), 185 (1983)
34. Knutson, B.: Facial expressions of emotion influence interpersonal trait inferences. *Journal of Nonverbal Behavior* **20**(3), 165–182 (1996)
35. Kong, Y., Jia, Y., Fu, Y.: Learning human interaction by interactive phrases. In: *European Conference on Computer Vision*, pp. 300–313 (2012)
36. Kostinger, M., Wohlhart, P., Roth, P., Bischof, H.: Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In: *IEEE International Conference on Computer Vision Workshop*, pp. 2144–2151 (2011)
37. Krizhevsky, A., Sutskever, I., Hinton, G.E.: ImageNet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems* (2012)
38. Kumar, N., Belhumeur, P., Nayar, S.: Facetracer: A search engine for large collections of images with faces. In: *European Conference on Computer Vision*, pp. 340–353. Springer (2008)
39. Lan, T., Sigal, L., Mori, G.: Social roles in hierarchical models for human activity recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2012)
40. Lee, D.H.: Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In: *International Conference on Machine Learning Workshop*, vol. 3, p. 2 (2013)
41. Levi, G., Hassner, T.: Emotion recognition in the wild via convolutional neural networks and mapped binary patterns. In: *ACM International Conference on Multimodal Interaction*, pp. 503–510 (2015)
42. Li, H., Lin, Z., Shen, X., Brandt, J., Hua, G.: A convolutional neural network cascade for face detection. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
43. Liu, M., Li, S., Shan, S., Chen, X.: AU-inspired deep networks for facial expression feature learning. *Neurocomputing* **159**, 126–136 (2015)
44. Liu, M., Li, S., Shan, S., Wang, R., Chen, X.: Deeply learning deformable facial action parts model for dynamic expression analysis. In: *Asian Conference on Computer Vision* (2014)
45. Liu, M., Shan, S., Wang, R., Chen, X.: Learning expressionlets on spatio-temporal manifold for dynamic facial expression recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2014)
46. Liu, P., Han, S., Meng, Z., Tong, Y.: Facial expression recognition via a boosted deep belief network. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1805–1812 (2014)
47. Liu, S., Yang, J., Huang, C., Yang, M.H.: Multi-objective convolutional learning for face labeling. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
48. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: *IEEE International Conference on Computer Vision* (2015)
49. Lucey, P., Cohn, J.F., Kanade, T., Saragih, J., Ambadar, Z., Matthews, I.: The extended cohn-kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression. In: *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 94–101 (2010)
50. Luo, P., Wang, X., Tang, X.: Hierarchical face parsing via deep learning. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2012)
51. Lyons, M.J., Budynek, J., Akamatsu, S.: Automatic classification of single facial images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **21**(12), 1357–1362 (1999)
52. Mollahosseini, A., Chan, D., Mahoor, M.H.: Going deeper in facial expression recognition using deep neural networks. In: *IEEE Winter Conference on Applications of Computer Vision* (2016)
53. Moody, J., Hanson, S., Krogh, A., Hertz, J.A.: A simple weight decay can improve generalization. *Advances in Neural Information Processing Systems* **4**, 950–957 (1995)
54. Ng, H.W., Nguyen, V.D., Vonikakis, V., Winkler, S.: Deep learning for emotion recognition on small datasets using transfer learning. In: *ACM International Conference on Multimodal Interaction*, pp. 443–449 (2015)
55. Opitz, M., Waltner, G., Poier, G., Possegger, H., Bischof, H.: Grid loss: Detecting occluded faces. In: *European Conference on Computer Vision* (2016)
56. Pantic, M., Cowie, R., D’Errico, F., Heylen, D., Mehu, M., Pelachaud, C., Poggi, I., Schroeder, M., Vinciarelli, A.: Social signal processing: the research agenda. In: *Visual analysis of humans*, pp. 511–538. Springer (2011)
57. Pantic, M., Valstar, M., Rademaker, R., Maat, L.: Web-based database for facial expression analysis. In: *IEEE International Conference on Multimedia and Expo* (2005)
58. Parkhi, O.M., Vedaldi, A., Zisserman, A.: Deep face recognition. In: *British Machine Vision Conference* (2015)
59. Pearson, K.: Note on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London* **58**, 240–242 (1895)
60. Pentland, A.: Social signal processing. *IEEE Signal Processing Magazine* **24**(4), 108 (2007)
61. Raducanu, B., Gatica-Perez, D.: Inferring competitive role patterns in reality TV show through nonverbal analysis. *Multimedia Tools and Applications* **56**(1), 207–226 (2012)
62. Ramanathan, V., Yao, B., Fei-Fei, L.: Social role discovery in human events. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2475–2482 (2013)
63. Ricci, E., Varadarajan, J., Subramanian, R., Rota Bulò, S., Ahuja, N., Lanz, O.: Uncovering interactions and interactors: Joint estimation of head, body orientation and f-formations from surveillance videos. In: *IEEE International Conference on Computer Vision* (2015)
64. Ruiz, A., Van de Weijer, J., Binefa, X.: From emotions to action units with hidden and semi-hidden-task learning. In: *IEEE International Conference on Computer Vision*, pp. 3703–3711 (2015)
65. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* **115**(3), 211–252 (2015)
66. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
67. Shan, C., Gong, S., McOwan, P.W.: Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing* **27**(6) (2009)
68. Sun, Y., Wang, X., Tang, X.: Sparsifying neural network connections for face recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
69. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
70. Tian, Y., Kanade, T., Cohn, J.F.: Facial expression recognition. In: *Handbook of face recognition*. Springer (2011)
71. Tian, Y.L., Kanade, T., Cohn, J.F.: Recognizing action units for facial expression analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **23**(2), 97–115 (2001)

72. Tianqi Chen Mu Li, Y.L.M.L.N.W.M.W.T.X.B.X.C.Z., Zhang, Z.: Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. In: NIPS Workshop on Machine Learning Systems (2016)
73. Trigeorgis, G., Snape, P., Nicolaou, M.A., Antonakos, E., Zafeiriou, S.: Mnemonic descent method: A recurrent process applied for end-to-end face alignment. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
74. Valstar, M.F., Mehu, M., Jiang, B., Pantic, M., Scherer, K.: Meta-analysis of the first facial expression recognition challenge. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics* **42**(4), 966–979 (2012)
75. Vinciarelli, A., Pantic, M., Bourlard, H.: Social signal processing: Survey of an emerging domain. *Image and Vision Computing* **27**(12), 1743–1759 (2009)
76. Vinciarelli, A., Pantic, M., Heylen, D., Pelachaud, C., Poggi, I., D’Errico, F., Schröder, M.: Bridging the gap between social animal and unsocial machine: A survey of social signal processing. *IEEE Transactions on Affective Computing* **3**(1), 69–87 (2012)
77. Wang, G., Gallagher, A., Luo, J., Forsyth, D.: Seeing people in social context: Recognizing people and social relationships. In: European Conference on Computer Vision, pp. 169–182 (2010)
78. Wang, J., Cheng, Y., Feris, R.S.: Walk and learn: Facial attribute representation learning from egocentric video and contextual data. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
79. Weng, C.Y., Chu, W.T., Wu, J.L.: RoleNet: Movie analysis from the perspective of social networks. *IEEE Transactions on Multimedia* **11**(2), 256–271 (2009)
80. Wu, Y., Ji, Q.: Constrained joint cascade regression framework for simultaneous facial action unit recognition and facial landmark detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
81. Yang, B., Yan, J., Lei, Z., Li, S.Z.: Aggregate channel features for multi-view face detection. In: International Joint Conference on Biometrics (2014)
82. Yang, B., Yan, J., Lei, Z., Li, S.Z.: Convolutional channel features. In: IEEE International Conference on Computer Vision (2015)
83. Yang, H., Zhou, J.T., Cai, J.: Improving multi-label learning with missing labels by structured semantic correlations. In: European Conference on Computer Vision, pp. 835–851 (2016)
84. Yang, S., Luo, P., Loy, C.C., Tang, X.: From facial parts responses to face detection: A deep learning approach. In: IEEE International Conference on Computer Vision (2015)
85. Yang, S., Luo, P., Loy, C.C., Tang, X.: Wider face: A face detection benchmark. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
86. Yao, A., Shao, J., Ma, N., Chen, Y.: Capturing au-aware facial features and their latent relations for emotion recognition in the wild. In: ACM International Conference on Multimodal Interaction, pp. 451–458 (2015)
87. Yu, H.F., Jain, P., Kar, P., Dhillon, I.: Large-scale multi-label learning with missing labels. In: International Conference on Machine Learning, pp. 593–601 (2014)
88. Yu, Z., Zhang, C.: Image based static facial expression recognition with multiple deep network learning. In: ACM International Conference on Multimodal Interaction, pp. 435–442 (2015)
89. Zafeiriou, S., Papaioannou, A., Kotsia, I., Nicolaou, M.A., Zhao, G.: Facial affect in-the-wild: A survey and a new database. In: IEEE Conference on Computer Vision and Pattern Recognition Workshop (2016)
90. Zelnik-Manor, L., Perona, P.: Self-tuning spectral clustering. In: NIPS, pp. 1601–1608 (2004)
91. Zhang, N., Paluri, M., Ranzato, M., Darrell, T., Bourdev, L.: Panda: Pose aligned networks for deep attribute modeling. In: IEEE Conference on Computer Vision and Pattern Recognition (2014)
92. Zhang, Z., Luo, P., Loy, C.C., Tang, X.: Learning deep representation for face alignment with auxiliary attributes. In: IEEE Transactions on Pattern Analysis and Machine Intelligence (2015)
93. Zhang, Z., Luo, P., Loy, C.C., Tang, X.: Learning social relation traits from face images. In: IEEE International Conference on Computer Vision (2015)
94. Zhang, Z., Luo, P., Loy, C.C., Tang, X.: Joint face representation adaptation and clustering in videos. In: European Conference on Computer Vision (2016)
95. Zhao, G., Huang, X., Taini, M., Li, S.Z., Pietikäinen, M.: Facial expression recognition from near-infrared videos. *Image and Vision Computing* **29**(9), 607–619 (2011)
96. Zhao, G., Pietikäinen, M.: Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **29**(6), 915–928 (2007)
97. Zhao, X., Liang, X., Liu, L., Li, T., Vasconcelos, N., Yan, S.: Peak-piloted deep network for facial expression recognition. In: European Conference on Computer Vision (2016)
98. Zhong, L., Liu, Q., Yang, P., Liu, B., Huang, J., Metaxas, D.N.: Learning active facial patches for expression analysis. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 2562–2569 (2012)
99. Zhu, X., Lei, Z., Liu, X., Shi, H., Li, S.Z.: Face alignment across large poses: A 3d solution. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)