

# A Unified Probabilistic Framework for Spontaneous Facial Action Modeling and Understanding

Yan Tong

GE Global Research Center

Niskayuna, NY12309

Jixu Chen and Qiang Ji

Rensselaer Polytechnic Institute

Troy, NY12180

## Abstract

**Facial expression is a natural and powerful means of human communication. Recognizing spontaneous facial actions, however, is very challenging due to subtle facial deformation, frequent head movements, and ambiguous and uncertain facial motion measurements. Because of these challenges, current research in facial expression recognition is limited to posed expressions and often in frontal view.**

A spontaneous facial expression is characterized by rigid head movements and nonrigid facial muscular movements. More importantly, it is the coherent and consistent spatial-temporal interactions among rigid and nonrigid facial motions that produce a meaningful facial expression. Recognizing this fact, we introduce a unified probabilistic facial action model based on the dynamic Bayesian network (DBN) to simultaneously and coherently represent rigid and nonrigid facial motions, their spatial-temporal dependencies, and their image measurements. Advanced machine learning methods are introduced to learn the model based on both training data and subjective prior knowledge. Given the model and the measurements of facial motions, facial action recognition is accomplished through probabilistic inference by systemically integrating visual measurements with the facial action model. **Experiments show that compared to the state-of-the-art techniques, the proposed system yields significant improvements in recognizing both rigid and nonrigid facial motions, especially for spontaneous facial expressions.**

## Index Terms

Facial Action Unit Recognition, Face Pose Estimation, Facial Action Analysis, Facial Action Coding System, Bayesian Networks.

### I. INTRODUCTION

Facial action is one of the most important sources of information for understanding emotional state and intention [1]. Spontaneous facial behavior is characterized by rigid head movement, nonrigid facial muscular movements, and their interactions. Rigid head movement characterizes the overall 3D head pose including rotation and translation. Nonrigid facial muscular movement results from the contraction of facial muscles and characterizes the local facial action at a finer level. The Facial Action Coding System (FACS) developed by Ekman and Friesen [2] is the most commonly used system for facial behavior analysis. Based on FACS, nonrigid facial muscular movement can be described by 44 facial action units (AUs), each of which is anatomically related to the contraction of a specific set of facial muscles.

An objective and noninvasive system for facial action understanding has applications in human behavior science, human-computer interaction, security, interactive games, computer-based learning, entertainment, telecommunication, and psychiatry. However, developing such a system faces several challenges:

- First, facial actions are rich and complex. Thousands of distinct nonrigid facial muscular movements (different AU combinations) have been observed so far [3], and most of them differ subtly in a few facial features.
- Second, **compared to the highly controlled conditions of posed facial expressions, spontaneous facial expressions often co-occur with natural head movement when people communicate with others. For example, a person may express his agreement by nodding his head and smiling simultaneously.** Hence, faces sometimes are partially occluded in some images. This makes it more challenging for accurately measuring facial motions.
- Third, most of the spontaneous facial expressions are activated without significant facial appearance changes, that is the amplitudes of the spontaneous facial expressions are smaller than those of the posed facial expressions. In addition, the spontaneous facial expression often has a slower onset phase and a slower offset phase compared to the posed facial

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

expression.

- Fourth, the spontaneous facial expression may have multiple apexes, and the expression does not always follow a neutral-expression-neutral temporal pattern [4] as for the posed facial expressions. Moreover, multiple facial expressions often occur sequentially.
- Fifth, the subtle facial deformations and frequent head movements make it more difficult to label the facial expression data. Hence, human labeling is difficult and less reliable.

Due to the low intensity, nonadditive effect, and individual difference of the spontaneous facial action as well as the image uncertainty, individually recognizing AUs is not accurate and reliable. Hence, understanding spontaneous facial action requires not only improving facial motion measurements, but more importantly, requires exploiting the spatial-temporal interactions among facial motions since it is these coherent, coordinated, and synchronized interactions of facial motions that produce a meaningful facial display. **By explicitly modeling and employing the spatiotemporal relationships in a facial action, the impact of these inaccurate or even erroneous facial motion measurements on the facial action recognition accuracy can be minimized. In addition, even if some facial motion measurements are missing due to occlusion, they can be inferred through their associations with other facial motions. Thus, the performance of facial action recognition can be improved.**

Our previous research on facial action unit recognition [5] shows that AU recognition benefits from explicitly modeling the relationships among AUs. However, the work is limited to AU recognition on nearly frontal view faces, and ignores the impact of head movement on AU measurements. Furthermore, the previous method focuses on modeling the semantic spatial relationships among AUs. In addition, the previous work recognizes AUs from posed facial expressions. This research, therefore, differs from our previous research in both theory and applications. Theoretically, this research models both the spatial and temporal interactions among AUs as well as modeling the interactions between the rigid motion (head pose) and the nonrigid motions (AUs). In application, this research focuses on recognizing spontaneous facial expressions, which are different from posed expressions. Moreover, recognizing spontaneous facial actions is much more challenging than recognizing posed facial actions.

In this paper, we introduce a probabilistic facial action model based on the dynamic Bayesian network (DBN) to simultaneously and coherently represent rigid head movement, nonrigid facial muscular movements, their spatial-temporal interactions, and their image observations in a spon-

taneous facial behavior. Advanced learning techniques are employed to construct the framework from both subjective knowledge and training data. Given the facial action model, facial action recognition is accomplished through probabilistic inference by systemically integrating the facial motion measurements with the facial action model.

The proposed facial action recognition system consists of two main stages: offline facial action model construction, and online facial motion measurement and inference. Specifically, using training data and subjective domain knowledge, the facial action model is constructed offline. During online recognition, as shown in Figure 1, various computer vision techniques are used to obtain measurements of both rigid (head pose) and nonrigid facial motions (AUs). These measurements are then used as evidences by the facial action model for inferring the true states of the rigid head pose and the nonrigid AUs simultaneously. **Currently, we only model left-right head movement since this type of head movement affects the AU measurement the most significantly compared to up-down and in-plane rotation and it appears frequently in spontaneous expressions. However, the system can be generalized to model the full range of head movement without changing the structure of the proposed facial activity model. The experiments show that compared to the state-of-the-art techniques, the facial action recognition with the proposed system is improved significantly, especially for spontaneous facial expressions.**

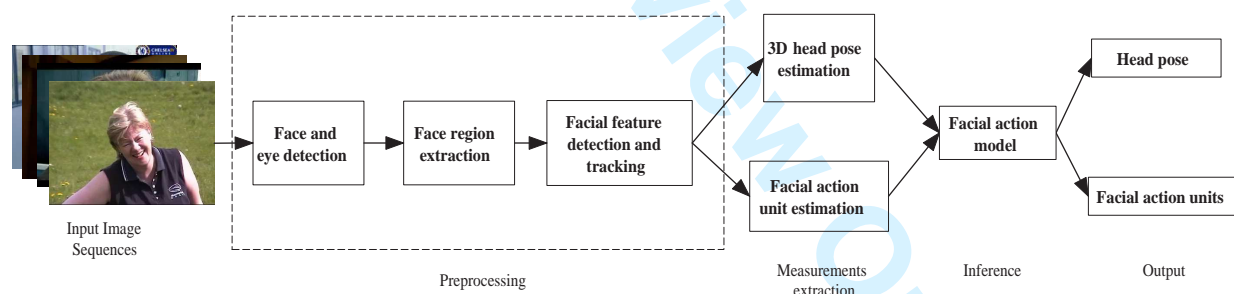


Fig. 1. The flowchart of the online facial action recognition system.

## II. RELATED WORK ON FACIAL ACTION RECOGNITION

Over the past fifteen years, there has been extensive research in computer vision on recognizing facial actions. Detailed surveys of previous work can be found in [6][7][8][9][1][10]. However, most of the previous research is limited to either estimating head poses on neutral faces or recognizing facial expressions/facial AUs on nearly frontal view faces due to the challenges

1  
2  
3  
4 mentioned previously. Since in spontaneous facial action people rarely express emotions without  
5 head movement, the assumption of nearly frontal view is not realistic.  
6

7  
8 Most recently, the research has been focused on improving facial action recognition under  
9 realistic conditions including spontaneous facial expressions with varying head pose. Recent work  
10 on facial action analysis can be classified into three groups based on how these methods deal with  
11 the relationships between the head movement and nonrigid facial muscular movements. The first  
12 group of methods [11][12] explicitly represents and recognizes the facial expression on 3D facial  
13 expression databases. The second group [13][14][15][16][17][18] assumes that the 3D head pose  
14 is independent of nonrigid facial muscular movements and estimates the 3D pose and nonrigid  
15 facial motions sequentially and separately. The third group [19][20][21][22][23] explicitly models  
16 the coupling between rigid head movement and nonrigid facial motions for facial action analysis,  
17 and simultaneously recovers the rigid and nonrigid facial motions. In contrast to recognizing  
18 facial action from deliberate facial display, automatic understanding of spontaneous facial action  
19 has been of greater interest in recent years [24][25][26][27][14][28][29][30][31][32]. In this  
20 section, we present a brief review of previous approaches to facial action analysis.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

### 31 *A. 3D Facial Expression Recognition*

32

33 Since nonrigid facial muscular movements produce 3D shape deformation of the facial surface,  
34 the first group of methods [11][12] extracts the 3D facial geometrical deformation caused by  
35 facial expression changes and recognizes the facial expression from facial deformations extracted  
36 from a 3D image database. Wang et al. [11] propose a 3D primitive surface feature based on  
37 principal curvature analysis on a 3D triangularized facial mesh. Facial expression recognition  
38 is performed on the static 3D images, whereas the dynamic nature of facial action is ignored.  
39 Instead of employing only the static 3D images, Chang et al. [12] use 3D facial geometry  
40 deformation for facial expression recognition from 3D videos.  
41  
42  
43  
44  
45  
46  
47

48 3D-based facial expression recognition is not affected by illumination changes and is inde-  
49 pendent of pose variation. However, these methods require high quality capture of texture and  
50 3D geometrical data, and thus are not applicable to many applications due to excessive hardware  
51 requirement. Furthermore, the computational complexity of the 3D methods on the dense data  
52 is considerable.  
53  
54  
55  
56  
57  
58  
59  
60

### B. Sequential and Separate Extraction of Rigid and Nonrigid Facial Motions

Assuming that 3D head pose is independent of nonrigid facial muscular movements, the second group of methods [13][14][15][16][17][18] estimates 3D pose and nonrigid facial motions sequentially and separately in two steps. Usually, a tracking process is performed first, and 3D pose is estimated from tracked salient facial feature points. Next, the facial expression is recognized from the pose-free facial texture or from the extracted nonrigid facial motions after eliminating the effect of pose.

Since rigid motion and nonrigid motions are nonlinearly coupled in the projected 2D facial shape/appearance, head pose estimation is not reliable under varying facial expressions. Likewise, because it is difficult to isolate the motion caused by facial expression from that caused by head movement, facial expression recognition is not accurate under varying head pose.

### C. Simultaneous Recovery of Rigid and Nonrigid Facial Motions

The coupling between rigid head movement and nonrigid facial motions can be modeled explicitly. The methods described in [19][20][21][23] estimate rigid head movement and nonrigid facial muscular movements simultaneously based on a 3D face model. The interaction between 3D head pose and facial expression is explicitly modeled as a nonlinear function.

Marks et al. [23] model the image sequence as a stochastic process generated by object motion including both rigid and nonrigid motion, object texture, and background texture. They track object motion, object texture, and background texture simultaneously by probabilistic filtering. Although they recover the 3D head pose and the nonrigid facial deformation, they do not perform facial expression or facial action recognition.

Vasilescu and Terzopoulos [22] propose a tensor face model based on multilinear image analysis, where the face image is generated from several modes such as identity of subject, head pose, and facial expression. The multilinear image analysis assumes that each mode is allowed to vary in turn, while the remaining factors are fixed. Hence, the head pose parameters and the facial expressions are explicitly modeled by two modes in the tensor representation, respectively. The core tensor represents the interaction between the modes.

**Based on an Active Appearance Model, Lucey et al [32] separate the rigid and non-rigid motions by dividing the shape parameters into “similarity parameters” characterizing the global head movement and “object-specific parameters” representing the non-rigid**

1  
2  
3  
4 **facial transformations. The AUs are recognized from a “similarity normalized” shape or**  
5 **appearance. However, it requires a specific AAM for each subject, which is not realistic**  
6 **in many applications.**  
7  
8

9  
10 Although the above methods successfully decouple the rigid and nonrigid facial motions,  
11 facial expression and head movement are recognized independently from the recovered rigid  
12 and nonrigid motions separately. These approaches neglect the interactions between rigid and  
13 nonrigid motions. Their measurements are, therefore, less robust, especially for spontaneous  
14 facial expressions.  
15  
16  
17

#### 18 19 *D. Facial Action Recognition from Spontaneous Facial Display*

20  
21 Most of the previous approaches recognize facial action from posed or deliberate facial  
22 displays. They are of limited practical use since only the spontaneous facial display can reflect  
23 the “true” emotion [27]. Moreover, posed facial display differs from spontaneous facial display  
24 in many aspects such as the magnitude, dynamic properties, and the interactions with head  
25 movement since they are initiated by two distinct areas of the brain [1]. Therefore, understand-  
26 ing spontaneous facial display is desirable and important for many real-world circumstances.  
27 Hence, facial action analysis recently began to focus on facial action recognition from spon-  
28 taneous facial behaviors rather than from posed facial behaviors. These approaches include  
29 recognizing spontaneous facial emotions/expressions [27][31][30] and recognizing facial AUs  
30 [24][25][26][14][28][29][32].  
31  
32  
33  
34  
35  
36  
37  
38

39 Unfortunately, besides performing the facial action recognition on different facial expression  
40 data (spontaneous facial expression data vs. posed facial expression data), most of the existing  
41 methods for recognizing spontaneous facial expressions employ the same techniques as for posed  
42 facial action, without exploiting the specific properties of spontaneous facial action. Therefore,  
43 these methods have the following limitations:  
44  
45  
46  
47

48 First, some of the current methods recognize each AU individually. However, since spontaneous  
49 facial action often produces subtle facial appearance changes rather than exaggerated appearance  
50 changes, recognizing AUs at low intensity levels is extremely difficult for current machine  
51 vision techniques. Therefore, recognition accuracy on the spontaneous facial display degrades  
52 significantly compared to accuracy on the posed facial display. In addition, for spontaneous facial  
53 expressions, AUs often occur in combination, and the appearance of an AU in a combination  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4 may be different from its appearance when occurring alone. This nonadditive effect makes it  
5  
6 more difficult to recognize AUs individually.

7  
8 Second, most of the current methods ignore the dynamic characteristics of AUs, which include  
9  
10 the dynamic development of each AU and the dynamic relationships among AUs. However,  
11  
12 recent psychological study [33][34][35] shows that the dynamic characteristics are crucial to  
13  
14 interpreting naturalistic human behavior. Valstar et al. [28] perform AU recognition for the  
15  
16 eyebrow movement, and find that spontaneous eyebrow motion can be explicitly distinguished  
17  
18 from posed eyebrow motion by employing the dynamic properties of the related AUs such as the  
19  
20 activating speed, magnitude, and the occurrence orders of AUs. Moreover, the research by Cohn  
21  
22 and Schmidt [26] shows that spontaneous smile usually has a relatively slower and smoother  
23  
24 onset, and that the intensity of lip corner motion is a strong linear function of time in contrast  
25  
26 to the posed smile.

27  
28 Finally, in spontaneous facial displays, facial expression changes are often accompanied with  
29  
30 natural head movements. Understanding spontaneous facial action should, therefore, deal with  
31  
32 the large facial shape/appearance variations caused by both head movement and nonrigid facial  
33  
34 muscular movements. Although most of the existing methods try to separate the facial motions  
35  
36 caused by facial muscular movements from those caused by head movement either manually  
37  
38 [24][14] or automatically [25][27][30][31], **they generally ignore the interactions between**  
39  
40 **the head movement and facial muscular movements, which is crucial for interpreting**  
41  
42 **spontaneous facial expressions. For example, if an up-down head movement is accompanied**  
43  
44 **by a smile facial expression, we may guess that the person wants to show an agreement.**

45  
46 In summary, current work focuses on either recognizing one type of facial motion while  
47  
48 ignoring the other, or recognizing both motions separately while ignoring their interactions.  
49  
50 Hence, these approaches cannot recognize facial actions reliably and robustly, especially for  
51  
52 spontaneous facial expressions. Current methods for spontaneous facial expression understanding  
53  
54 use the same techniques as those for posed expression recognition. Few existing methods consider  
55  
56 the spatial-temporal interactions among facial motions. Although our previous work [36] models  
57  
58 the interactions between rigid and nonrigid motions, its modeling is limited to mainly spatial  
59  
60 interactions, and it is limited to recognizing the facial action from posed facial displays. In  
contrast, this work explicitly models the spatial-temporal interactions among rigid and nonrigid  
facial motions and uses the model to improve the facial expression recognition.



### III. FACIAL ACTION MODELING

#### A. Overview of the Facial Action Model

A spontaneous facial action consists of rigid head motion, nonrigid facial deformations, and their interactions. In the scenario of facial action analysis from 2D images, the 2D facial shape can be generated by a stochastic process from three hidden causes we want to infer: head pose, 3D facial shape, and nonrigid facial muscular movements. The 3D facial shape characterizes an intrinsic property of a subject, and it differs from subject to subject. The nonrigid facial muscular movements, which are systematically represented by a set of AUs, cause the 3D shape deformation of the facial surface. The head pose characterizes the overall head movement including rotation and translation, and causes the changes in the position and shape of the 2D face on the image. In addition, through various computer vision techniques, we can obtain measurements for these hidden causes. Figure 2(a) shows a graphical model representing such causal relationships among different elements of a facial action.

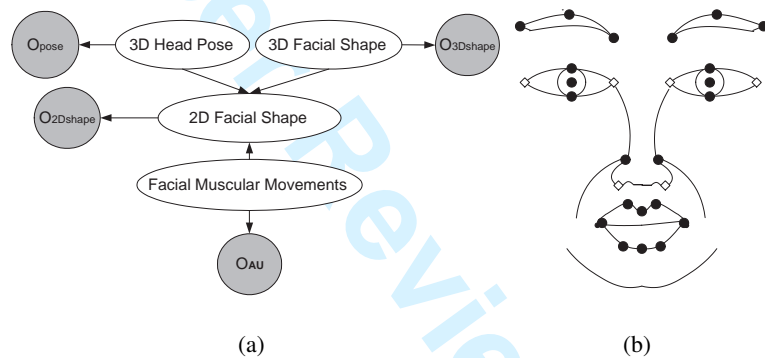


Fig. 2. (a) A graphical model to represent the causal relationships among elements of a facial action, where the shaded nodes represent the measurements of the connected hidden nodes. (b) Facial feature points on a frontal view face: the black dots represent the local shape points, whereas the white diamonds represent the global shape points.

Based on the causal relationships shown in Figure 2(a), we propose to use a Bayesian network (BN) to model 3D facial shape, facial muscular movements, 3D head pose, 2D facial shape, and their relationships. A BN is a directed acyclic graph (DAG) where each node represents a random variable, and the link between two variables characterizes the causal relationship between them. Such a model is capable of representing the statistical dependencies among the rigid motion, nonrigid motions, and their interactions through the 2D facial shape. Furthermore, the nodes in a BN can be grouped into hidden nodes and measurement nodes. The 3D facial shape, facial muscular movements, 3D head pose, and 2D facial shapes are modeled as hidden nodes, and their true states are what we want to infer from their measurements through the model. And so, we

associate each hidden node with a measurement node (shaded) representing the observation of the corresponding hidden node through some computer vision techniques, as shown in Figure 2(a).

Given the model, facial action recognition is to find the optimal states of rigid motion (head pose) and nonrigid facial muscular movements (AUs) by maximizing the joint probability of pose and AUs given their measurements as follows:

$$pose^*, \mathbf{AU}^* = \underset{pose, \mathbf{AU}}{\operatorname{argmax}} p(pose, \mathbf{AU} | O_{pose}, O_{\mathbf{AU}}, O_{3Dshape}, O_{2Dshape}) \quad (1)$$

where  $\mathbf{AU}$  is the set of all AUs of interest,  $O_{pose}$ ,  $O_{\mathbf{AU}}$ ,  $O_{3Dshape}$ , and  $O_{2Dshape}$  denote the measurements of the head pose, AUs, 3D facial shape, and 2D facial shape, respectively, as shown in Figure 2(a).

In the next several sections, we gradually show how the relationships in Figure 2(a) can be expanded and enriched, based on which we can solve for Eq. (1).

### B. Modeling Rigid Motion with 2D Global Shape

In this research, the shape of a 3D face can be represented by a vector of 28 facial feature points as shown in Figure 2(b). The 28 points are located around each facial component (e.g. mouth, eye, nose, etc.). The 28 points are further divided into global feature points as represented by the white diamonds (e.g., the points on the lower nose corners and the points at eye corners) and local feature points (e.g., the points on the eye lids and the points on the lips) as represented by the black dots.

Given a 3D face, the deformation of a 2D facial shape reflects the action of both head pose and facial muscular movements. Specifically, the head pose and facial muscular movements may affect different sets of facial feature points. The global facial feature points are relatively invariant to the facial muscular movements, and their movements are primarily caused by head pose. For instance, whether the eye is open or closed does not change the positions of the eye corners. On the other hand, the local facial feature points are not only affected by the head pose; they are also sensitive to the facial muscular movements. **Based on this observation, the 2D facial shape can be decomposed into a global facial shape and local facial component shapes, as shown in Figure 2(b). They form a two-level hierarchical structure, as shown in Figure 3(b).**

The 2D global shape  $\mathbf{S}_g$  is the projection of the 3D global feature points on the image plane; therefore, it is directly affected by the 3D facial shape and the head pose. The 3D facial shape

governs the shape of the 2D global shape, whereas the 3D head pose controls both the position and shape of the 2D global shape. This causal dependency can be represented by a directed link from the head pose/3D facial shape to the 2D global shape  $S_g$ , as shown in Figure 3(a).

### C. Modeling the Relationship between 2D Global Shape and 2D Local Facial Component Shapes

The 2D local facial shape is partitioned into four components: eyebrow, eyes, nose, and mouth. The two eyes (or eyebrows) are considered as one facial component because of their symmetry. Each 2D local facial component shape is indirectly affected by the rigid head movement through the 2D global shape  $S_g$ . Given the 2D global shape  $S_g$ , the position (center) of each local facial component can be roughly estimated, independent of the head pose. For example, the center of eye can be determined, given the eye corners, which are parts of the global shape. Hence, this causal relationship can be represented by a directed link from the 2D global shape to each 2D local facial component shape as illustrated in Figure 3(b).

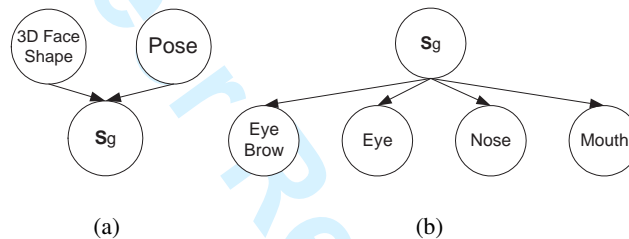


Fig. 3. (a) The head pose and 3D facial shape directly affect the 2D global shape  $S_g$ . (b) The causal relationship between the 2D global shape  $S_g$  and each 2D local facial component shape.

### D. Modeling the Relationships between Nonrigid Motion and 2D Local Facial Component Shapes

The nonrigid facial muscular movements produce significant changes in the 3D shape of the facial component. These 3D facial muscular movements can be systematically represented by AUs, which are anatomically related to the contraction of the facial muscles as defined in [37]. For example, activating AU27 (mouth stretch) will produce a widely open mouth; and activating AU4 (brow lowerer) makes the eyebrows lower and pushed together. In this work, we intend to recognize a set of commonly occurring AUs<sup>1</sup>.

Since the 3D shape of each facial component is determined by the related AUs, the 2D local facial component shape is also controlled by the AUs, besides the rigid head movement.

<sup>1</sup>AU1 (Inner brow raiser), AU2 (Outer brow raiser), AU4 (Brow lowerer), AU5 (Upper lid raiser), AU6 (Cheek raiser and lid compressor), AU7 (Lid tightener), AU9 (Nose wrinkler), AU12 (Lip corner puller), AU15 (Lip corner depressor), AU17 (Chin raiser), AU23 (Lip tightener), AU24 (Lip Presser), AU25 (Lips part), and AU27 (Mouth stretch).

For instance, there are six AUs (AU12, AU17, AU23, AU24, AU25, and AU27) controlling mouth movements, and three AUs (AU1, AU2, and AU4) controlling eyebrow movements. We, therefore, could directly connect the related AUs to the corresponding 2D local facial components to represent the causal relationships among them. For example, AU1 (Inner brow raiser), AU2 (Outer brow raiser), AU4 (Brow lowerer) could be connected to the 2D eyebrow node, while AU9 (nose wrinkler) could be connected to the 2D nose node, since activating AU9 will pull the infraorbital triangle upwards.

However, directly connecting all related AUs to one facial component would result in too many AU combinations, most of which rarely occur in daily life. For example, based on the analysis of the training data, there are only 8 common AU or AU combinations for the mouth in spite of 64 potential AU combinations. Thus, only a set of common AU or AU combinations, which produce significant nonrigid facial actions, is sufficient to control the shape variations of the facial component. As a result, a set of intermediate nodes (i.e., “ $C_B$ ”, “ $C_E$ ”, and “ $C_M$ ” for eyebrow, eye, and mouth, respectively) are explicitly introduced to model the correlations among AUs and to reduce the number of AU combinations. For example, the intermediate node “ $C_M$ ” has 8 states, each of which represents a common AU or AU combination controlling mouth movement. Figure 4(a) shows the modeling of the relationships between the nonrigid 3D facial motions (AUs) and the 2D local facial component shapes.

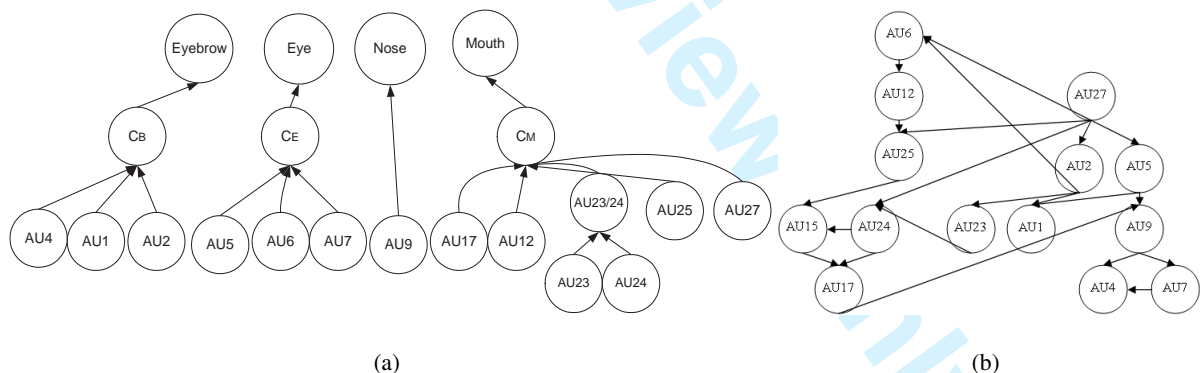


Fig. 4. (a) The relationship between the AUs and the local facial component shape. Intermediate nodes ( $C_B$ ,  $C_E$ , and  $C_M$ ) are introduced to model the correlations among AUs and to reduce the number of AU combinations to model. (b) The semantic relationships among AUs, which are necessary to produce a coherent and meaningful display. Details may be found in [5].

### E. Semantic AU Relationships Modeling

So far, we have modeled relationships between rigid head motion and nonrigid facial motions through the 2D facial shapes, but have not discussed the modeling of relationships among nonrigid facial motions, i.e., modeling spatial dependencies among facial action units. Based

on our previous study in [5], AUs are spatially and semantically related in order to create a meaningful facial display. In particular, there are two important semantic relationships among the AUs: co-occurrence relationships and mutually exclusive relationships. The co-occurrence relationships characterize some groups of AUs, which usually appear together to show meaningful facial displays. For example, if the mouth and the eyes are observed to be widely open, the eyebrows are most likely raised up, since it implies a surprise expression. On the other hand, based on the alternative rules provided in the FACS manual, some AUs are mutually exclusive, since “it may not be possible anatomically to do both AUs simultaneously” or “the logic of FACS precludes the scoring of both AUs” [37]. For instance, the lips cannot be parted as AU25 (lips part) and pressed as AU24 (lip presser) simultaneously. These semantic relationships among AUs are especially important for understanding spontaneous facial display. For example, the enjoyment or Duchenne smiles are often accompanied by the facial muscular movements related to AU12 (lip corner puller) and AU6 (cheek raiser), whereas the miserable or dampened smiles are often produced by AU12 (lip corner puller) and AU15 (lip corner depressor)[38][39]. Figure 4(b) shows a BN that models that spatial dependencies among some AUs. Details about this model and its acquisition may be found in [5].

#### IV. MODELING THE DYNAMIC RELATIONSHIPS AMONG AUs

In a spontaneous facial action, AUs are not only spatially dependent on each other, but also show strong temporal dependencies to represent different naturalistic facial expressions. Nishio et al. [40] have shown that when the mouth moves prior to the eye movement, a smile expression is mostly interpreted as a smile of enjoyment. On the contrary, when the eyes move prior to the mouth movement, it is mostly interpreted as a dampened smile.

Generally speaking, there are two types of temporal dependencies among AUs: intra-dependency and inter-dependency. Intra-dependency characterizes the self temporal evolution of an AU, while inter-dependency captures temporal dependencies between different AUs, i.e., an AU will be activated following the activation of another AU. For example, in a spontaneous smile, AU12 (lip corner puller) is usually first activated to express a slight emotion; then, with the increasing of emotion intensity, AU6 (cheek raiser) is activated in an average of 0.4 second after the activation of AU12 [39]; and after both the actions reach their apexes simultaneously, AU6 is relaxed and AU12 is gradually released before both of them return to the neutral state. Furthermore, due to the

variability among individuals and different contexts, the dynamic relationships among AUs are stochastic. Systemically capturing such temporal dependencies among AUs and incorporating them into the facial action recognition process is especially important for interpretation of spontaneous facial behaviors.

#### A. A DBN for Modeling Dynamic Dependencies among AUs

In this paper, we propose to use a Dynamic Bayesian Network (DBN) to model and learn the dynamic dependencies among AUs. A DBN is a directed acyclic graphical model, which models the temporal evolution of a set of random variables  $\mathbf{X}$  over time [41]. It represents a generalization of the traditional dynamic models including the Hidden Markov Models and Kalman Filtering. Let  $\mathbf{X}^t$  represent a set of random variables at a discrete time slice  $t$ . A DBN is defined as  $B = (G, \Theta)$ , where  $G$  is the model structure, and  $\Theta$  represents the model parameters, i.e., the conditional probability tables (CPTs) for all nodes. There are two assumptions in the DBN model: first, we assume that the system is first-order Markovian, i.e.,  $P(\mathbf{X}^{t+1}|\mathbf{X}^0, \dots, \mathbf{X}^t) = P(\mathbf{X}^{t+1}|\mathbf{X}^t)$ ; and second, it is assumed that the process is stationary, i.e., that the transition probability  $P(\mathbf{X}^{t+1}|\mathbf{X}^t)$  is the same for all  $t$ . Therefore, a DBN  $B$  can be also defined by a pair  $(B_0, B_{\rightarrow})$ : (1) the static network  $B_0 = (G_0, \Theta_0)$ , as shown in Figure 5(a), captures the static distribution over all variables  $\mathbf{X}^0$ ; and (2) the transition network  $B_{\rightarrow} = (G_{\rightarrow}, \Theta_{\rightarrow})$ , as shown in Figure 5(b), specifies the transition probability  $P(\mathbf{X}^{t+1}|\mathbf{X}^t)$  for all  $t$  in a finite time slices  $T$ .

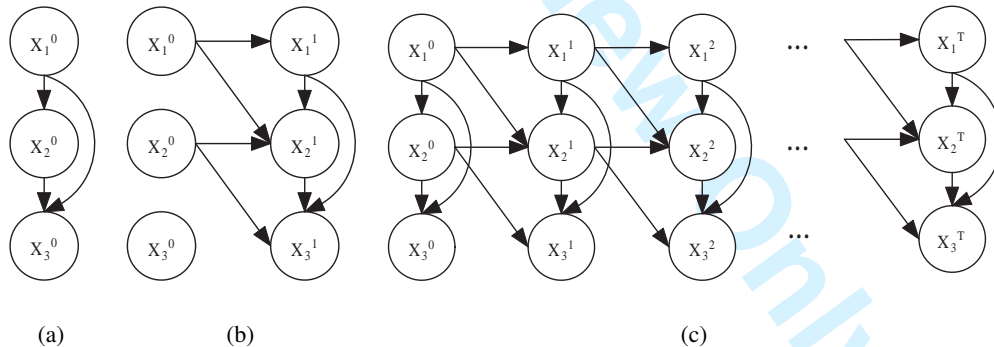


Fig. 5. A pair of (a) static network  $B_0$  and (b) transition network  $B_{\rightarrow}$  defines the dynamic dependencies for three random variables  $X_1, X_2, X_3$ . (c) The corresponding “unrolled” DBN for  $T + 1$  time slices.

Given a DBN model, the joint probability over all variables  $\mathbf{X}^0, \dots, \mathbf{X}^T$  can be factorized by “unrolling” the DBN into an extended static BN, as shown in Figure 5(c), whose joint probability is computed as follows:

$$P(\mathbf{x}^0, \dots, \mathbf{x}^T) = P_{B_0}(\mathbf{x}^0) \prod_{t=0}^{T-1} P_{B_{\rightarrow}}(\mathbf{x}^{t+1}|\mathbf{x}^t), \quad (2)$$

where  $\mathbf{x}^t$  represents the sets of values taken by the random variables  $\mathbf{X}$  at time  $t$ ,  $P_{B_0}(\mathbf{x}^0)$  captures the joint probability of all variables in the static BN  $B_0$ , and  $P_{B_{\rightarrow}}(\mathbf{x}^{t+1}|\mathbf{x}^t)$  represents the transition probability and can be decomposed as follows based on the conditional independencies encoded in the DBN:

$$P_{B_{\rightarrow}}(\mathbf{x}^{t+1}|\mathbf{x}^t) = \prod_{i=1}^N P_{B_{\rightarrow}}(x_i^{t+1}|pa(X_i^{t+1})), \quad (3)$$

where  $pa(X_i^{t+1})$  represents the parent configuration of variable  $X_i^{t+1}$  in the transition network  $B_{\rightarrow}$ , and  $N$  represents the number of random variables in  $\mathbf{X}^t$ . Hereafter,  $pa^j(X)$  represents the  $j^{th}$  parent configuration of variable  $X$  in a given network structure.

In this paper, we intend to use DBN to learn and model two types of temporal relationships for AUs at two adjacent time slices. The intra-dependency is modeled as an arc linking an  $AU_i$  node at time  $t - 1$  ( $AU_i^{t-1}$ ) to that at time  $t$  ( $AU_i^t$ ) and depicts how a single AU develops over time. The inter-dependency is modeled as an arc from  $AU_i$  node at time  $t - 1$  to  $AU_j$  at time  $t$  and represents the pairwise dynamic dependency between two different AUs.

However, we should notice that the temporal relationships are critically depending on the temporal resolution of the image sequences, i.e., the frame rate of the recorded videos. The learned temporal relationships may not be valid on the image sequences collected under a different frame rate. For example, the image sequences in Cohn-Kanade database [42] are recorded with 12 fps, whereas those in ISL facial expression database [43] are collected with 30 fps. Therefore, a time duration of 1/6 second consists of two frames for Cohn-Kanade database [42] and five frames for ISL database [43]. Therefore, we consider the temporal relationships between two fixed-size time durations instead of the temporal relationships between two successive frames. In addition, for each AU, its temporal evolution consists of a complete temporal segment lasting from 1/4 of a second (e.g., a blink) to several minutes (e.g., a jaw clench) as described in [4]. If we choose a single frame as one time duration, we may capture many irrelevant events; whereas if we choose many frames as a duration, the dynamic relationships may not be captured. Hence, based on an analysis of the databases we use as well as on the temporal characteristics of the AUs we intend to recognize, the time duration is empirically set as 1/6 second in this work. This way, if  $AU_i$  appears within the  $t^{th}$  time duration, it is counted as “presence” at  $t$ , i.e.,  $AU_i^t = 1$ . Hereafter, we use “slice” instead of “duration” for a generalized representation in the later presentation.

### B. Constructing the Initial DBN

In this work, each AU is represented by a binary value  $[0, 1]$  representing its absence/presence status. An AU is in the “presence” status, when it is activated and is at one of the three temporal states (onset, apex, and offset); whereas it is in the “absence” status, when it is at the neutral state. Since the “presence/absence” of an AU at time  $t$  depends not only on its state in previous time slice, but also on the states of other AUs,  $P(AU_j^t | AU_j^{t-1}, AU_i^{t-1})$  is used to capture the dynamic relationships between  $AU_i$  and  $AU_j$  as well as the dynamic evolution of  $AU_j$  itself. For example, the positive dependency between two AUs in adjacent slices is computed as follows:

$$P(AU_j^t = 1 | AU_j^{t-1} = 1, AU_i^{t-1} = 1) = \frac{N_{AU_j^t + AU_j^{t-1} + AU_i^{t-1}}}{N_{AU_j^{t-1} + AU_i^{t-1}}}, \quad (4)$$

where  $N_{AU_j^{t-1} + AU_i^{t-1}}$  is the total number of the events that the AU combination  $(AU_j + AU_i)$  is present in the  $(t-1)^{th}$  slice, regardless of the presence of other AUs, and  $N_{AU_j^t + AU_j^{t-1} + AU_i^{t-1}}$  is the total number of the events that  $AU_j$  is present in the  $t^{th}$  slice and the AU combination  $(AU_j + AU_i)$  is present in the  $(t-1)^{th}$  slice in the databases. The other probabilities are computed similarly.

For initialization, **the intra-dependency and inter-dependency are partially learned from two posed facial expression databases:** namely, Cohn-Kanade facial expression database [42] and ISL facial expression database [43]. Furthermore, in order to recognize AUs in spontaneous facial expression, we will later refine the initial dynamic relationships among AUs in the spontaneous facial expression using another database containing natural facial expressions such as [44] and [45]. Based on the data analysis from the two databases [42][43], we can find that the statistical information, extracted from training data, is consistent with the dynamic relationships among AUs found in the psychological studies [39]. For example,  $P(AU_6^t = 1 | AU_6^{t-1} = 0, AU_{12}^{t-1} = 1) = 0.1 > P(AU_6^t = 1 | AU_6^{t-1} = 0, AU_{12}^{t-1} = 0) = 0.02$ , which means AU6 (cheek raiser) occurs mostly after AU12 (lip corner puller) is activated.  $P(AU_5^t = 1 | AU_5^{t-1} = 0, AU_2^{t-1} = 1) = 0.15 > P(AU_5^t = 1 | AU_5^{t-1} = 0, AU_2^{t-1} = 0) = 0.02$ , which means AU5 (lid raiser) is activated mostly after AU2 (outer brow raiser) is activated.

If the probability  $P(AU_j^t = 1 | AU_j^{t-1} = 0, AU_i^{t-1} = 1)$  is higher than a predefined threshold  $T_{up}$  or the probability  $P(AU_j^t = 1 | AU_j^{t-1} = 1, AU_i^{t-1} = 0)$  is lower than a predefined threshold  $T_{bottom}$ , we assume that there is a strong dynamic dependency between  $AU_i$  and  $AU_j$ , which can be modeled with an inter-slice link from  $AU_i^{t-1}$  to  $AU_j^t$  in the DBN. For example, the link from



$AU_{12}^{t-1}$  to  $AU_6^t$  represents the dynamic dependency between AU6 (cheek raiser) and AU12 (lip corner puller), since AU6 is present mostly after AU12 is activated. This dynamic dependency can be also derived from the psychological studies [39]. The link from  $AU_2^{t-1}$  to  $AU_5^t$  means that AU5 (lid raiser) is activated mostly after AU2 (outer brow raiser) is activated. This way, using the statistics extracted from the two databases, an initial transition network is manually constructed as in Figure 6(a).

### C. Learning DBN Model

Given a set of observed data  $D = \{D_1, \dots, D_M\}$ , we can refine the initial DBN model with a structure learning algorithm, i.e., finding a DBN structure  $G$  that best fits the observed data. For learning a DBN model, the training data  $D$  should be divided into  $S$  sequences, each of which contains an “expressive” segment of *neutral-expression1*– $\dots$ –*expressionN*–*neutral* with length  $M_s$  so that  $\sum_s M_s = M$ , where  $M$  is the total number of training images. As mentioned above, a DBN consists of two parts ( $B_0$  and  $B_{\rightarrow}$ ), therefore, we should learn both of them from the training data. To evaluate the fitness of the network, we need to define a scoring function. The score of a DBN model can be defined based on the Bayesian Information Criterion (BIC) [46], i.e.,

$$Score(B) = \log P(B, D) = \log P(B) + \log P(D|B), \quad (5)$$

where  $\log P(B)$  is the log prior probability of the network structure, and  $\log P(D|B)$  is the log likelihood, which can be computed approximately as follows:

$$\log P(D|B) \approx \log P(D|G, \widehat{\Theta}_G) - \frac{\log M}{2} Dim_G, \quad (6)$$

where the first term evaluates how well the network  $B$  fits the data  $D$ , the second term is a penalty relating to the complexity of the network,  $\widehat{\Theta}_G$  is the set of parameters of  $G$  which maximizes the likelihood of the data,  $M$  is the number of training data, and  $Dim_G$  is the number of parameters.

Instead of giving an equal prior probability  $P(B)$  to all possible structures, we assign a higher probability to the manually constructed initial structure  $B_{init}$  shown in Figure 6a as in [47]. Given the scoring function, the structure learning is performed by searching the network with the highest score among the possible network structures.

Specifically, the score for a DBN model can be decomposed into two parts as follows:

$$Score(B) = Score_{B_0} + Score_{B_{\rightarrow}}, \quad (7)$$

where  $Score_{B_0}$  and  $Score_{B_{\rightarrow}}$  represent the score of the static network and the score of the transition network, respectively. Consequently, we can learn the structure of  $B_0$  and the structure of  $B_{\rightarrow}$  separately.

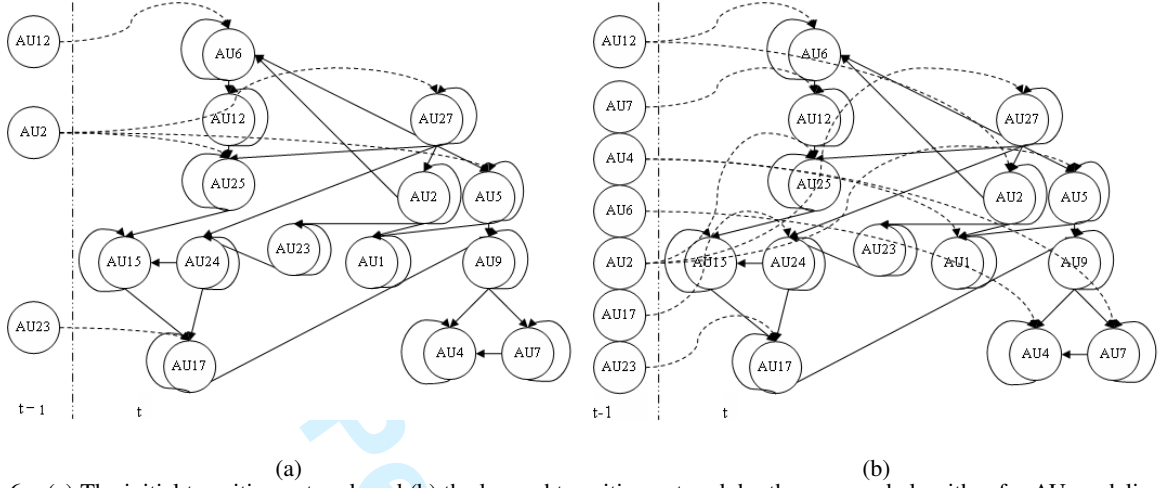


Fig. 6. (a) The initial transition network and (b) the learned transition network by the proposed algorithm for AU modeling. The self-arrow at each AU node indicates the temporal relationship of a single AU from the previous time slice to the current time slice. The dashed line with arrow from  $AU_i$  at time  $t - 1$  to  $AU_j$  ( $j \neq i$ ) at time  $t$  indicates the pairwise dynamic dependency between different AUs.

1) *Learning the static network:* The static network  $B_0$  models the semantic spatial relationships among AUs within a time slice as discussed in section III-E. Following Eq.(5) and Eq.(6), the score for the static network is defined as follows:

$$Score_{B_0} = \log P(B_0) + \sum_i \sum_j \sum_k N_{i,j,k}^0 \log \hat{\theta}_{i,j,k}^0 - \frac{\log M}{2} Dim_{B_0}, \quad (8)$$

where  $\theta_{i,j,k}^0$  represents the parameters for the static network  $B_0$ , i.e.,  $\theta_{i,j,k}^0 = P(X_i^0 = k | pa^j(X_i^0))$  with  $i$  ranges over all variables,  $j$  ranges over all parent configuration of one specific variable  $X_i^0$ , and  $k$  ranges over all states of  $X_i^0$  in  $B_0$ .  $N_{i,j,k}^0 = \sum_M I(X_i^0 = k, pa^j(X_i^0))$ , where  $I(\cdot)$  is an indicator function so that  $I(\cdot) = 1$  if the argument is true, otherwise  $I(\cdot) = 0$ .

Given the definition of  $Score_{B_0}$ , we employ an iterated hill climbing algorithm [48] to search the optimal network structure. First, an initial network structure  $B_{0init}$  is manually constructed by combining the data analysis from the two databases [42][43] and the domain knowledge from the FACS rules [37]. Then, starting from  $B_0^0 = B_{0init}$ ,  $Score_{B_0}$  is computed as in Eq. (8) for each nearest neighbor of  $B_0^0$ , which is generated from  $B_0^0$  by adding, deleting, or reversing a single arc that is subject to the acyclicity constraint and the limitation on the upper bound of parent nodes. In this way, the network structure that has the maximum score among all of the nearest neighbors is selected as the static network  $B_0$  as shown in Figure 4(b). Details on

learning the static BN structure for modeling the semantic AU relationships may be found in our previous work [5].

2) *Learning the transition network:* Learning the transition network is more complicated than learning the static network. The transition network  $B_{\rightarrow}$  consists of two types of links: inter-slice links and intra-slice links. The inter-slice links are the dynamic links connecting the temporal variables of two successive time slices. In contrast, the intra-slice links connect the variables within a single time slice, which are same as the static network structure. The score of the transition network is defined as follows:

$$Score_{B_{\rightarrow}} = \log(P(B_{\rightarrow})) + \sum_i \sum_j \sum_k N_{i,j,k}^{\rightarrow} \log \hat{\theta}_{i,j,k}^{\rightarrow} - \frac{\log(M-S)}{2} Dim_{B_{\rightarrow}}, \quad (9)$$

where  $M - S$  is the total number of pair-wise transitions between two successive slices in the training data, and  $\theta_{i,j,k}^{\rightarrow}$  represents the parameters for the transition network  $B_{\rightarrow}$  and is defined as  $\theta_{i,j,k}^{\rightarrow} = P(X_i^t = k | pa^j(X_i^t))$  for the node  $X_i^t$  at  $k^{th}$  state given its  $j^{th}$  parent configuration in  $B_{\rightarrow}$ .  $N_{i,j,k}^{\rightarrow}$  accounts for the number of the instances of transitions and is defined as  $N_{i,j,k}^{\rightarrow} = \sum_s \sum_t I(X_i^t = k, pa^j(X_i^t))$  where  $I(\cdot)$  is an indicator function.

Given the definition of a score for the transition network as in Eq. (9), we need to identify a transition network structure with the highest score by a searching algorithm subject to some coherent constraints on the transition network. First, the variables  $\mathbf{X}^0$ , as shown in Figure 5(b) do not have parents. Second, the inter-slice links can only have one direction, that is from the previous time slice to the current time slice. Finally, based on the stationary assumption, both the inter-slice links and intra-slice links should be repeated for the time slice  $t \in [1, T]$ . **Furthermore, since the multi-wise relationships are hard to capture and vary significantly from people to people, for better generalization, we only capture the strong pairwise dynamic dependencies among AUs that are true for most people and ignore the weak temporal relationships that are mostly person-dependent.** Therefore, an additional constraint is imposed so that each node  $X_i^{t+1}$  has at most two parents from the previous time slice  $t$ . We then apply the same hill climbing technique [48] to identify the transition network structure.

Figure 6(b) shows the learned transition network by the proposed learning algorithm. Compared with the manually constructed initial transition network in Figure 6a, the learned structure better reflects the dynamic relationships among AUs in the training data. For example, the dynamic link from  $AU_4^{t-1}$  to  $AU_7^t$  means that the eye lids intend to be narrowed by activating

AU7 (lid tightener) with the increasing intensity of AU4 (brow lowerer). And the dynamic link from  $AU_{17}^{t-1}$  to  $AU_{24}^t$  means that before the lips are pressed together (AU24), it is most likely the chin boss<sup>2</sup> is already moved upward by activating AU17 (chin raiser).

## V. FACIAL MOTION MEASUREMENTS AND THEIR MODELING

In the facial action model so far, the head pose, AUs, 3D facial shape, 2D global facial shape, 2D local facial component shapes, and the intermediate nodes are hidden nodes. The hidden nodes can be divided into “observable” hidden nodes and “non-observable” hidden nodes. The “observable” hidden nodes are those nodes we can acquire their measurements using some computer vision techniques. They include those nodes representing head pose, AUs, 3D facial shape, 2D global facial shape, and 2D local facial component shapes. The “non-observable” hidden nodes are the remaining hidden nodes, i.e. the intermediate nodes.

To acquire the measurements for the “observable” hidden nodes, we first perform face and eye detection on neutral face with frontal view by using a boosted eye detection algorithm based on recursive nonparametric discriminant analysis (RNDA)[49]. Once the face and eye centers are detected, the face region is normalized and is convolved pixel by pixel by a set of multiscale and multiorientation Gabor filters. Next, the 28 facial feature points, as shown in Figure 2(b), are detected on the neutral and frontal face similar to [50]. After that, we obtain the measurement of the 3D facial shape by personalizing a trained generic 3D shape model. **Specifically, the x- and y- coordinates of the generic 3D shape model are adapted to current subject based on the detected positions of facial feature points on the frontal view face. Due to the unknown depth information, the z-coordinates of the generic 3D shape model are scaled based on the size of the same frontal view face to approximate the face depth for each individual.**

Afterward, the measurements of 2D global shape and local facial component shapes are obtained by tracking the 28 facial feature points in each image frame. **Specifically, we use active shape models and Gabor wavelet for facial feature tracking as described in [51].**

Furthermore, based on the personalized 3D facial shape and the tracked global facial feature points, **three face pose angles (i.e., pan, tilt and roll) are estimated through the weak perspective projection model by using a technique as described in [20].** Since the in-plane

<sup>2</sup>Chin boss is a term that defined in Table 1-1 (page 3) of the FACS manual[37]. It means “the skin covering the bone of the chin”.

1  
2  
3  
4 rotation and translation can be compensated by image normalization given the eye position, we  
5 only focus on modeling and recognizing the head pose variation from the left-right rotation.  
6 Furthermore, for simplicity, the continuous pan angle is discretized into frontal, left, and right  
7 face pose measurement.  
8  
9

10  
11 Given the normalized face image, **we also extract the measurement for each AU through a**  
12 **general-purpose learning mechanism based on Gabor wavelet-based feature representation**  
13 **and AdaBoost classification following the work in [14]. Since we intend to recognize the**  
14 **AUs under varying face pose, for each AU, an AdaBoost classifier is constructed for each**  
15 **state of the head pose (frontal, left, and right), respectively. Assuming that the face pose**  
16 **varies smoothly over time, for each AU, the AdaBoost classifier, which corresponds to**  
17 **the face pose estimated in the previous frame, is used to extract its measurement for the**  
18 **current frame.**  
19  
20

21  
22 In this way, we obtain the image measurements for 3D facial shape, 2D global shape, 2D  
23 local component shapes, 3D head pose, and AUs. **To incorporate the measurements into**  
24 **the facial action model, a node (shaded) is introduced to represent each measurement, as**  
25 **shown in Figure 7.** In addition, a measurement link between a hidden node and its measurement  
26 is introduced to model the causal relationship between the hidden node and its measurement.  
27 The conditional probability that quantifies each measurement link may be used to model the  
28 measurement accuracy.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

## 39 VI. A COMPLETE FACIAL ACTION MODEL AND ITS PARAMETRIZATION

### 40 A. A Comprehensive Model for Facial Action Understanding

41  
42 Now we are ready to present the complete DBN model for facial action modeling , as shown  
43 in Figure 7. Specifically, there are three layers in the proposed model. The first layer consists of  
44 the 3D head pose, 3D facial shape, and the 2D global shape  $S_g$ , and it models the effect of rigid  
45 head motion and the 3D facial shape on the 2D global shape. The second layer contains a set  
46 of 2D local shapes corresponding to the facial components as well as their relationships to the  
47 global shape. The third layer includes a set of AUs and their spatial interactions and dynamic  
48 interactions through the dynamic links. The third layer also models the effect of AUs on the  
49 shapes of local facial components through the intermediate nodes.  
50  
51  
52  
53  
54  
55

56 We employ the first layer as the global constraint for the overall system so that it will guarantee  
57 globally meaningful facial action. Meanwhile, the local structural details of the facial components  
58  
59  
60

are constrained not only by the local shape parameters, but also are characterized by the nonrigid facial muscular movements (represented by the related AUs) through the interactions between the second layer and the third layer. In addition, the interactions between the rigid head motion and the nonrigid facial actions are indirectly modeled through the 2D global and local facial component shapes. Finally, the facial motion measurements are systematically incorporated into the model through the shaded nodes.

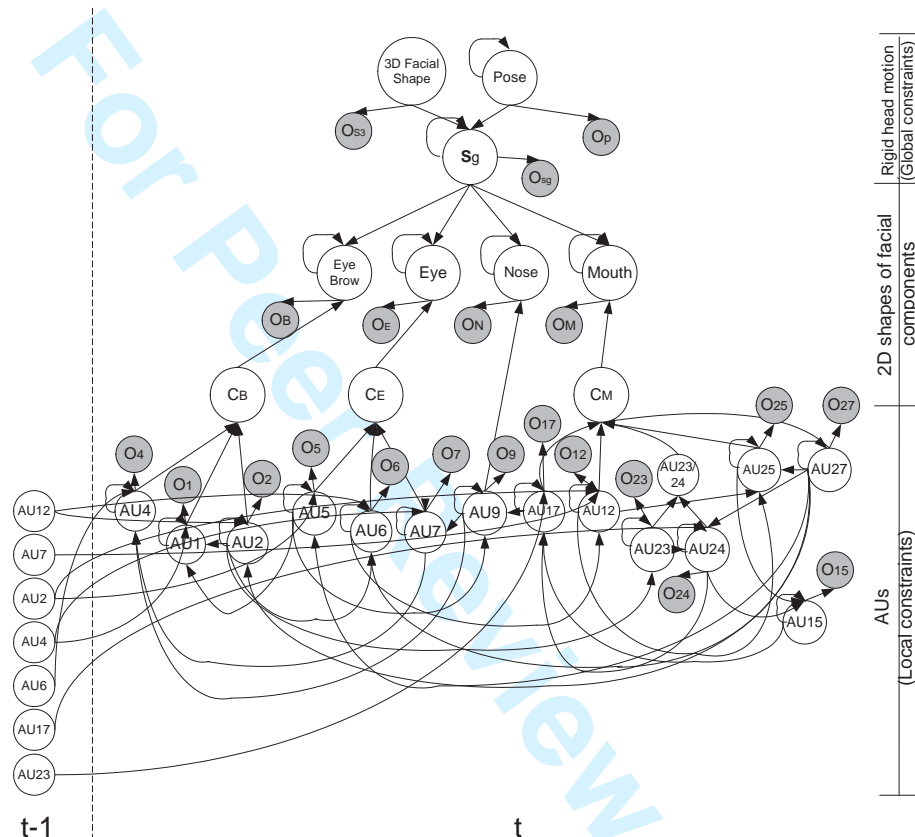


Fig. 7. The complete DBN model for facial action understanding. The shaded node indicates the observation for the connected hidden node. The self-arrow at the hidden node represents its temporal evolution from previous time frame to the current time frame. The link from  $AU_i$  at time  $t - 1$  to  $AU_j$  ( $j \neq i$ ) at time  $t$  indicates the dynamic dependency between different AUs.

This model, therefore, completely characterizes the spatial and temporal dependencies between rigid and nonrigid facial motions and accounts for the uncertainties in facial motion measurements.

### B. Model Learning And Parametrization

Given the model structure shown in Figure 7, we need to define the states for each node and, then, learn the model parameters associated with each node. For each node without parents, it is parameterized by its prior probability. For the continuous node with discrete/continuous

parents, it is characterized by Conditional Probabilistic Distribution (CPD) defined as conditional probability of a node  $X$ , given its parents  $pa(X)$ , i.e.,  $p(X|pa(X))$ ; whereas for the discrete node with discrete parents, it is characterized by the Conditional Probabilistic Table (CPT) defined as  $p(X|pa(X))$  similar to CPD.

Specifically, the head pose is represented by three different views: left, frontal, and right, which correspond to three discrete states ( $Pose \in 0, 1, 2$ ) in the proposed system. The prior information of the pose  $p(Pose)$  can be learned from the training images. The 3D facial shape  $S_{3D}$  is characterized by a continuous 3D shape vector consisting of 28 facial feature points from neutral faces; therefore, the node  $S_{3D}$  has continuous state. The 2D global shape  $\mathbf{S}_g$  is represented by a continuous shape vector consisting of global feature points, whereas the 2D local shape of the  $j^{th}$  facial component, such as the eye and the eyebrow, is represented by a continuous shape vector  $\mathbf{S}_{l_j}$  containing the corresponding local feature points. Hence, both the 2D global shape  $\mathbf{S}_g$  and 2D local facial component shape  $\mathbf{S}_{l_j}$  have continuous states. And each AU has two discrete states, which represent the “presence/absence” states of the AU.

Given the head pose at  $k^{th}$  state  $Pose = k$  and the 3D facial shape  $\mathbf{S}_{3D} = s_{3D}$ , the CPD of  $\mathbf{S}_g$  can be defined as we see here [52]:

$$p(\mathbf{S}_g = s_g | Pose = k, \mathbf{S}_{3D} = s_{3D}) = (2\pi)^{-\frac{d_g}{2}} |\Sigma_{gk}|^{-\frac{1}{2}} \exp\left(-\frac{\gamma_{gk}^2}{2}\right) \quad (10)$$

where  $d_g$  is the dimension of the 2D global shape  $\mathbf{S}_g$ , and  $\gamma_{gk}^2$  is defined as a Mahalanobis distance:

$$\gamma_{gk}^2 = (s_g - \mathbf{W}_{gk} * s_{3D} - \mu_{gk})^T \Sigma_{gk}^{-1} (s_g - \mathbf{W}_{gk} * s_{3D} - \mu_{gk}) \quad (11)$$

with the corresponding mean shape vector  $\mu_{gk}$ , regression matrix  $\mathbf{W}_{gk}$ , and covariance matrix  $\Sigma_{gk}$ . Based on the conditional independence embedded in the BN, we can learn the  $\mu_{gk}$ ,  $\mathbf{W}_{gk}$  and  $\Sigma_{gk}$  locally, as shown in Figure 3(a), from the training data consisting of 2D global shape, 3D facial shape, and head pose.

The intermediate nodes (i.e.,  $C_B$ ,  $C_E$ , and  $C_M$ ) are discrete nodes, each state of which represents a specific AU/AU combination related to a local facial component. For instance, “ $C_M$ ” has 8 states, each of which represents the presence of an AU or AU combination related to mouth movement. The CPT (conditional probabilistic table)  $p(C_i|pa(C_i))$  for each intermediate node is manually specified based on the data analysis. For example, we assign  $p(C_B = 0|AU1 =$

0,  $AU2 = 0, AU4 = 0$ ) = 0.9 if the eyebrow is at the neutral state, whereas  $p(C_B = 1|AU1 = 1, AU2 = 1, AU4 = 0)$  = 0.9 if the eyebrow is entirely raised up.

For each local shape component node (i.e., *EyeBrow*, *Eye*, *Nose*, and *Mouth*), its CPD is parameterized as a Gaussian distribution. For example, for the *EyeBrow* node, let  $\mathbf{S}_B$  denote the 2D local shape of eyebrow, the CPD of *EyeBrow*  $p(\mathbf{S}_B = s_B|\mathbf{S}_g = s_g, C_B = k)$  is assumed to satisfy a Gaussian distribution as follows:

$$p(\mathbf{S}_B = s_B|\mathbf{S}_g = s_g, C_B = k) = (2\pi)^{-\frac{d_B}{2}} |\Sigma_{Bk}|^{-\frac{1}{2}} \exp\left(-\frac{\gamma_{Bk}^2}{2}\right) \quad (12)$$

where  $d_B$  is the dimension of the 2D local shape  $\mathbf{S}_B$  for the eyebrow, and  $\gamma_{Bk}^2$  is defined as a Mahalanobis distance:

$$\gamma_{Bk}^2 = (s_B - \mathbf{W}_{Bk} * s_g - \mu_{Bk})^T \Sigma_{Bk}^{-1} (s_B - \mathbf{W}_{Bk} * s_g - \mu_{Bk}) \quad (13)$$

with corresponding mean shape vector  $\mu_{Bk}$ , regression matrix  $\mathbf{W}_{Bk}$ , and covariance matrix  $\Sigma_{Bk}$ . Given the training data of 2D local shape of eyebrow, 2D global shape, and the related AUs, we can learn the parameters  $\mu_{Bk}$ ,  $\mathbf{W}_{Bk}$ , and  $\Sigma_{Bk}$  locally. The parameters for *Eye*, *Nose*, and *Mouth* are defined and learned similarly to those of *EyeBrow*.

The CPD of each measurement node, given its parent is learned to reflect the measurement accuracy of the computer vision technique. For example,  $p(O_{AU_i}|AU_i)$  represents the measurement accuracy with the corresponding AdaBoost classifier.

The CPTs for the static links among all AUs are learned simultaneously in the local DBN model, as shown in the Figure 6(b). Finally, we learn the transition probability for the temporal links of the DBN. Specifically, based on Eq. (3), learning the transition probability associated with the temporal link is to learn the conditional probability of the temporal variable  $X_i$ , given its parent configuration in the transition network  $p_{B\rightarrow}(X_i^{t+1}|pa(X_i^{t+1}))$ . Since the training data is complete, we can learn the conditional probability  $p_{B\rightarrow}(X_i^{t+1}|pa(X_i^{t+1}))$  using Maximum Likelihood (ML) estimation method.

### C. Analysis of Training Data Size for Model Learning

As described previously, based on the conditional independencies encoded in the DBN model, the facial action model as shown in Figure 7 can be decomposed into a set of smaller local structures. As a result, the parameters of the facial action model can be learned individually in each local structure. Hence, fewer training data are required for



learning parameters in a smaller model than are required for a large network structure. To evaluate the quantity of training data needed for learning the facial action model, we perform a sensitivity study of model learning on different amount of training data. For this purpose, the Kullback-Leibler (KL) divergences of the parameters are computed versus the number of training samples. Specifically, two local models are learned: the dynamic AU model capturing the spatiotemporal relationships among all target AUs as shown in the third layer of Figure 7; and a local shape component model (“Eye Brow”) capturing the relationships among  $S_g$ , “Eye Brow”, “ $C_B$ ”, the related AUs, and their corresponding measurement nodes as shown in Figure 7.

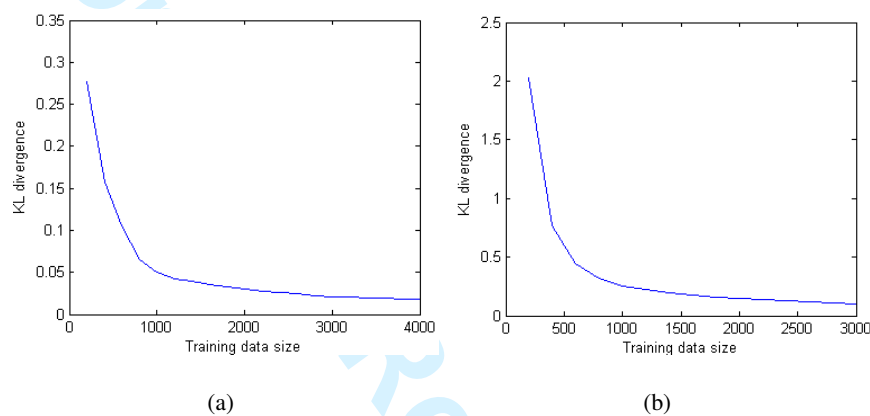


Fig. 8. The KL divergence of the parameters versus the training data size for (a) the dynamic AU model and (b) “EyeBrow” model, respectively.

The experimental results reported in Figure 8 show that for the dynamic AU model, the learning process requires a total of only 2000 training samples for all target AUs since all AUs have binary states. For the “Eye Brow” model, however, since the states of  $S_g$  and “Eye Brow” are continuous, the training process requires 1500 training samples. Since we have more than 2000 training samples available in the databases, the training data for the model learning is sufficient. Furthermore, our recent studies in [53][54] show that we can significantly reduce the training data by incorporating some qualitative knowledge in the model learning. By exploiting some qualitative constraints on the parameters, we can achieve the similar learning results with only one tenth of the training data required in the conventional learning techniques.

## VII. FACIAL ACTION INFERENCE

Once the measurement nodes are observed, we can infer the facial action by finding the most probable explanation (MPE) of the evidence, as shown in Eq.(1). The advantage of using MPE

is that it allows us to infer all the variables of interest simultaneously, instead of inferring each variable individually; therefore, it simultaneously finds the most probable state combination of head pose and AUs.

Denote  $O_{S_3}$ ,  $O_p$ ,  $O_{S_g}$ ,  $O_{S_{l_j}}$ , and  $O_{AU_i}$  as the measurements of the 3D face, head pose, the 2D global shape, the  $j^{th}$  2D local facial component shape, and  $AU_i$ , respectively. Based on the conditional independence encoded in the DBN, the inference can be factorized as follows:

$$p(pose, AU_{1...N} | O_{S_3}, O_p, O_{S_g}, O_{S_{l_1...M}}, O_{AU_{1...N}}) = c \int_{S_{3D}, S_g, S_{l_1...M}} \sum_{C_{1...K}} \{p(pose)p(S_{3D}) \quad (14)$$

$$p(O_{S_3} | S_{3D})p(S_g | S_{3D}, pose) \left[ \prod_j^M p(S_{l_j} | pa(S_{l_j})) \right] \left[ \prod_k^K p(C_k | pa(C_k)) \right] \left[ \prod_i^N p(AU_i | pa(AU_i)) \right]$$

$$p(O_{S_g} | S_g)p(O_p | pose) \left[ \prod_j^M p(O_{S_{l_j}} | S_{l_j}) \right] \left[ \prod_i^N p(O_{AU_i} | AU_i) \right]$$

where  $c$  is a normalization constant,  $M$  is the number of local facial components,  $N$  is the number of target AUs, and  $K$  is the number of the intermediate nodes. The factorized probabilities in Eq. (14) are the CPDs/CPTs that are learned as discussed in Section VI.

Therefore, the true joint states of head pose and the AUs can be inferred simultaneously, given the measurements of the 3D face, head pose, the 2D global shape, the 2D local shapes, and the AUs through probabilistic inference. Specifically, we employ the junction tree inference algorithm in the Bayes Net Toolbox by Murphy [55] to infer the true states of head pose and the AUs.

## VIII. EXPERIMENTAL RESULTS

### A. Facial Expression Databases

The proposed facial action analysis system is trained and tested on FACS labeled images from three databases. The first database is Cohn-Kanade DFAT-504 database [42], which consists of more than 100 subjects covering different races, ages, and genders. This database has been widely used for evaluating facial AU recognition system. However, the image sequences in Cohn-Kanade database only contain frontal view face images. In order to demonstrate our system under more natural and realistic circumstance, we also constructed our own database (ISL multiview facial expression database[43]), which consists of 40 image sequences from 8 subjects containing the target AUs. The ISL database <sup>3</sup> is collected under uncontrolled indoor illumination and

<sup>3</sup>More details about ISL multiview facial expression database can be found at <http://www.ecse.rpi.edu/homepages/cvrl/database/database.html>.

1  
2  
3  
4 background. The subjects are instructed to perform the target AUs or the basic facial expressions  
5 (e.g., smiling or surprising) while turning around their head. Hence, the face undergoes large  
6 face pose variations ( $-30^\circ$  to  $30^\circ$  from left to right) and significant facial expression changes  
7 simultaneously.  
8  
9

10  
11 However, both the Cohn-Kanade database and the ISL database contain posed facial expres-  
12 sions, whereas the interactions between the rigid motion and nonrigid facial muscular movements  
13 and the dynamic and semantic relationships among AUs in the spontaneous facial actions are  
14 different from those in the posed facial actions. Therefore, the relationships that are learned  
15 from the posed facial actions would bias the recognition against the spontaneous facial actions.  
16 Hence, the proposed system will be also trained and tested on spontaneous facial expression  
17 databases, whereas the training and testing procedures are the same as those on the posed  
18 facial expression databases. Therefore, we extend our work for recognizing spontaneous facial  
19 actions. Specifically, the proposed system is trained and tested on spontaneous facial expression  
20 databases which consists of image videos collected through three sources: (1) Multiple Aspects  
21 of Discourse (MAD) research lab at the University of Memphis [45]; (2) Belfast natural facial  
22 expression database [44]; and (3) videos obtained from the website (<http://www.youtube.com/>).  
23 In these image sequences, the subjects are displaying various spontaneous facial expressions  
24 with natural head movements.  
25  
26

27  
28 For this study, all the image sequences in the databases (Cohn-Kanade database, ISL database,  
29 and spontaneous facial expression databases) are coded into AUs frame by frame. For each AU,  
30 the positive samples are chosen as the images containing the target AU at different intensity  
31 levels, and the negative samples are selected as those images without the target AU regardless  
32 the presence of other AUs. For training the facial shape models, we also manually marked the  
33 28 feature points on some images from these databases. **In the following, we will perform  
34 experimental validation on the proposed system and compare the performance of the  
35 proposed system with the state-of-the-art techniques [14][5] on these image databases,  
36 respectively.**  
37

#### 38 *B. Evaluation on Cohn-Kanade DataBase*

39  
40 We first evaluate our system on the Cohn-Kanade database [42] for AU recognition to demon-  
41 strate the system performance on the standard database. The database is divided into eight  
42 sections, each of which contains images from different subjects. Each time, we use seven sections  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

for training and the remaining section for testing so that the training and testing set are mutually exclusive. The average recognition performance is computed on all sections.

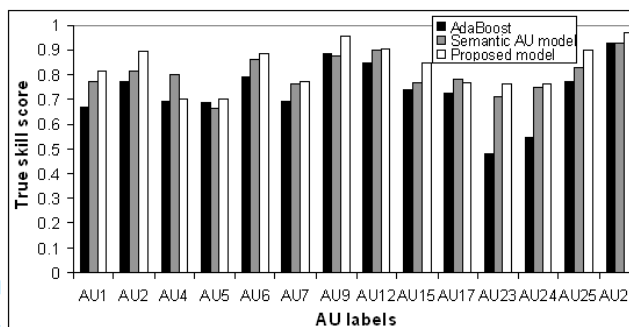


Fig. 9. Comparison of AU recognition results on the novel subjects in Cohn-Kanade database by using the AdaBoost classifier [14] (black bar), using the semantic AU model as in [5] (grey bar), and using the proposed model (white bar), respectively, based on the true skill score (Hansen Kuiper Discriminant), which is the difference between the correct-positive recognition rate and the false-positive rate.

Figure 9 shows the performance for generalization to novel subjects in Cohn-Kanade database of using the AdaBoost classifiers alone [14], using the semantic AU model that only models AU relationships as in [5], and using the proposed model, respectively. The AdaBoost classifiers [14] achieve an average correct-positive recognition rate of 80.6% and an average false-positive rate of 7.84% for the 14 target AUs. By employing the relationships among AUs, the semantic AU method increases the average correct-positive rate to 85.8% and reduces the false-positive rate to 5.54%. With the use of the proposed model, our system achieves an average correct-positive rate of 88.3% and reduced false-positive rate to 5.17%. Compared with the semantic AU model, the improvement by using the facial action model is not that significant. That is because the Cohn-Kanade database only contains the images with frontal view faces and posed facial expressions. We will further demonstrate the benefits of the proposed facial action model in the latter sections.

### C. Evaluation under Realistic Environment



Fig. 10. An example image sequence from ISL database where the subject is laughing with left-to-right head rotation.

In order to demonstrate the robustness of the proposed system, we perform experiments on our own database under realistic environment where the face undergoes facial expression and face pose changes simultaneously. Figure 10 shows an example of image sequence from the ISL database [43] where the subject is laughing and turning around his head.

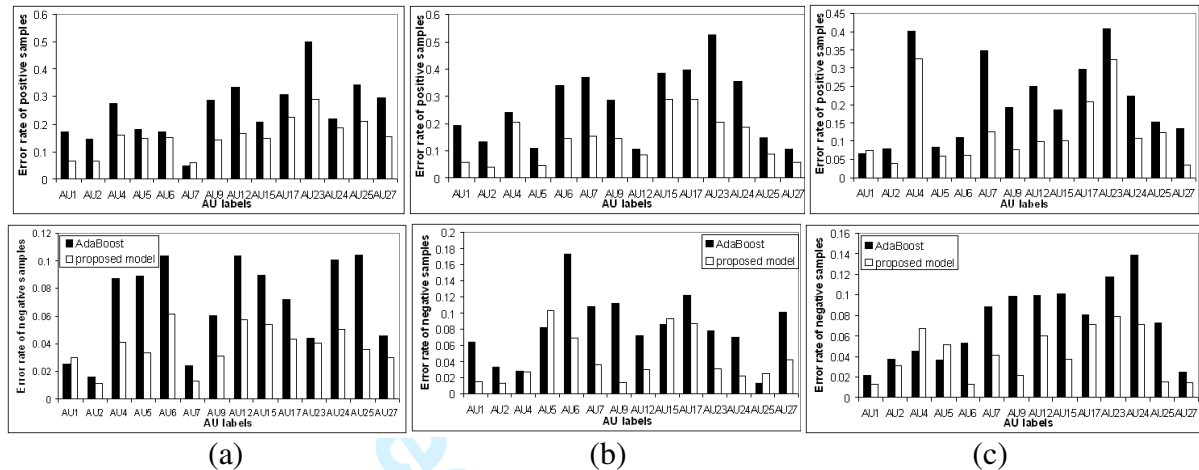


Fig. 11. AU recognition results under realistic circumstance for (a) frontal-view faces (left column), (b) left-view faces (middle column), and (c) right-view faces (right column). The first row demonstrates average error rate of positive samples. The second row displays average error rate of negative samples. In each figure, the black bar denotes the result by the AdaBoost classifier [14], and the white bar represents the result by using the proposed model.

The system is evaluated based on the leave-one-subject-out cross validation. The system performance is reported in Figure 11. Compared to the AU recognition by the AdaBoost classifiers [14] only, we can find that using the proposed facial action model: (1) for the frontal-view face, the average relative error rate of positive samples (positive error rate) decreases by 37.5%, and the average relative error rate of negative samples (negative error rate) decreases by 44.7%; (2) for the right-view face, the average relative positive error rate decreases by 40%, and the average relative negative error rate decreases by 42.2%; and (3) for the left-view face, the average relative positive error rate decreases by 46.1%, and the average relative negative error rate decreases by 46.8%. Here, the relative error rate is defined as the ratio of the error rate of the proposed method to the error rate of the AdaBoost method [14]. Especially for the AUs that are difficult to be recognized, the system performance is greatly improved. For example, for AU23 (lip tighten), its positive error rate decreases from 52.3% to 20.5%, and its negative error rate decreases from 7.7% to 3.1% for the left view face; the positive error rate of AU7 (lid tighten) decreases from 34.8% to 12.5% for the right view face; and the positive error rate of AU6 (cheek raiser and lid compressor) is increased from 34% to 14.5% with a significant drop of negative error rate

decreasing from 17.3% to 6.95% for the left view face.

We also perform pose estimation on the image sequences through the probabilistic inference. As shown in Table I, pose estimation by the proposed method is also improved compared to the pose measurement obtained by a method as in [20]. The improvement comes from modeling the interactions of head pose with AUs through the 2D facial shapes. As a result, the erroneous pose measurement is compensated by its relationships to the AUs. In summary, the proposed system significantly improves AU recognition and pose estimation simultaneously.

TABLE I

COMPARISON OF POSE ESTIMATION BY USING [20] AND USING THE DBN INFERENCE THROUGH THE PROPOSED MODEL.

view	[20]	proposed method
frontal	93%	94.3%
right	94.4%	96.1%
left	86.7%	93.4%

#### D. Evaluation on Spontaneous Facial Expression Database



Fig. 12. An example image sequence from Belfast database where the subject is talking with natural head movement.

Instead of recognizing posed facial actions, it is more important to recognize spontaneous facial actions. Therefore, in the third set of the experiments, the system is trained and tested on the spontaneous facial expression databases to demonstrate the system robustness for recognizing the spontaneous facial action. Currently, the combined spontaneous facial expression database contains 74 image sequences from 13 subjects, where 63 image sequences are used for training, and 11 image sequences for testing. Figure 12 shows an example of image sequence from the Belfast database [44] where the subject is talking with natural head movement.

Since there are major differences between the spontaneous facial action and posed facial action, as described in Section I, we should learn a new facial action model from the spontaneous training data to capture the relationships among AUs, the coupling between the rigid motion and nonrigid motion, and the dynamics in order to robustly recognize spontaneous facial action. In this work, we intend to recognize 12 target AUs, which frequently occur in the spontaneous facial expression database. Figure 13 shows the learned DBN model from the spontaneous facial expression database.

Figure 14 shows the average AU recognition performance on the spontaneous facial expression database of using the AdaBoost classifier alone [14], using the semantic AU model as in [5], and using the proposed facial action model respectively. The AdaBoost classifiers [14] achieve an average positive error rate of 44% and an average negative error rate of 8.58% for the 12 target AUs. By employing the relationships among AUs, the semantic AU model decreases the average positive error rate to 36.5% and decreases the negative error rate to 6.6%. With the use of the proposed complete facial action model, the average positive error rate further decreases to 24.3%, and the negative error rate decreases to 5.3%.

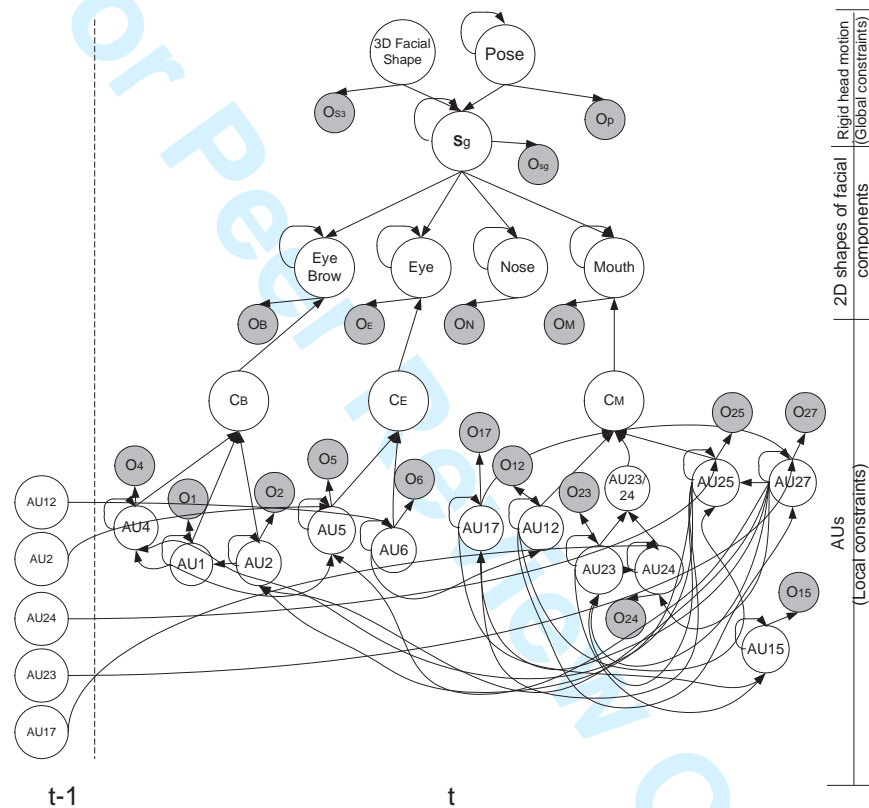


Fig. 13. Here we have the complete DBN model for spontaneous facial action recognition.

Compared to the recognition results from the AdaBoost classifiers [14], the system performance is greatly improved by using the proposed facial action model for some AUs. For example, the positive error rate of AU23 (lip tighten) decreases from 94.4% to 25.9% with a moderate increasing of negative error rate (from 3.6% to 5.8%); the positive error rate of AU12 (lip corner puller) decreases from 53% to 37.8% with a significant drop of negative error rate (from 32% to 12.8%); and the positive error rate of AU2 (outer brow raiser) decreases from 26.9% to 16.9% with a decreasing of negative error rate from 6.2% to 3.1%.

Furthermore, since the spontaneous facial action is often accompanied by natural head movements, only employing the relationships among AUs are not sufficient to deal with the facial appearance variations due to varying head pose. By incorporating the relationships between head pose and AUs through the 2D facial shapes, the complete facial action model further improves the system performance compared to the semantic AU model [5].

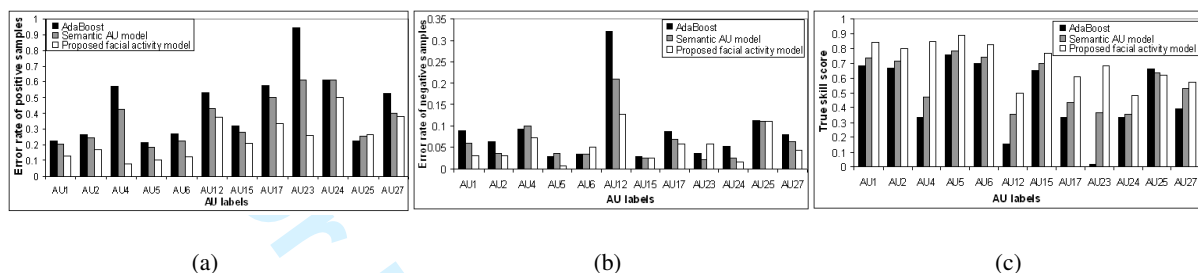


Fig. 14. AU recognition results on spontaneous facial expression database: (a) average positive error rates, (b) average negative error rates, and (c) average true skill scores (the difference between the correct-positive recognition rate and the false-positive rate). In each figure, the black bar denotes the result by using the AdaBoost classifier [14], the grey bar represents the result by using the semantic AU model [5], and the white bar represents the result by using the proposed facial action model.

## IX. CONCLUSION AND FUTURE WORK

The recent research shows that the spatiotemporal relationships among the rigid and non-rigid facial motions are important for spontaneous facial action analysis and understanding. Therefore, we propose a unified facial action model based on the DBN to systematically discover and learn such relationships, and then combine them with the image observations to perform a robust and reliable recognition of spontaneous facial action. The experiments show that compared to the state-of-the-art techniques [14][5], the proposed system yields significant improvements in both pose estimation and AU recognition, especially for the spontaneous facial expressions. The performance improvement is especially impressive for some difficult AUs such as AU6 (cheek raiser). The performance improvements come mainly from combining the facial action model with the facial measurements. Specifically, the erroneous AU measurements can be compensated by the model's built-in spatial and temporal relationships among AUs and the built-in relationships between rigid head motions and the nonrigid facial motions. The important lesson we can learn from this work is that for a robust visual interpretation and understanding, solely improving the computer vision techniques will not be enough. It is important to capture the prior knowledge or context in a probabilistic manner and systematically



combines the captured knowledge with the improved visual measurements to achieve a robust and accurate visual interpretation.

**Currently, the system can process about 2.1 frames/second in  $320 \times 240$  images on a 2.8GHz Pentium IV PC. The system could be sped up by optimization of the code and by implementing the DBN inference in C++ instead of in Matlab®.** The future work will focus on the following aspects. First, we plan to extend the work to include continuous measurement of head pose and to include both tilt and pan rotations. Second, we will extend the DBN model to systematically model human labeling errors by introducing another layer of nodes to represent the labeling confidence. For applications, we would like to apply the framework to distinguish faked facial expression from genuine and natural facial expressions. The proposed method may be also applied to perform “soft” biometrics for human identification based on their facial behaviors.

#### REFERENCES

- [1] M. Pantic and M. Bartlett, “Machine analysis of facial expressions,” in *Face Recognition*, K. Delac and M. Grgic, Eds. Vienna, Austria: I-Tech Education and Publishing, 2007, pp. 377–416.
- [2] P. Ekman and W. V. Friesen, *Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Palo Alto, CA: Consulting Psychologists Press, 1978.
- [3] K. Scherer and P. Ekman, *Handbook of Methods in Nonverbal Behavior Research*. Cambridge, UK: Cambridge University Press, 1982.
- [4] M. Pantic and I. Patras, “Dynamics of facial expression: recognition of facial actions and their temporal segments from face profile image sequences,” *IEEE Trans. SMC - PartB: Cybernetics*, vol. 36, no. 2, pp. 433–449, April 2006.
- [5] Y. Tong, W. Liao, and Q. Ji, “Facial action unit recognition by exploiting their dynamic and semantic relationships,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 29, no. 10, pp. 1683–1699, October 2007.
- [6] M. Pantic and L. J. M. Rothkrantz, “Automatic analysis of facial expressions: The state of the art,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 12, pp. 1424–1445, 2000.
- [7] B. Fasel and J. Luetin, “Automatic facial expression analysis: A survey,” *Pattern Recognition*, vol. 36, pp. 259–275, 2003.
- [8] R. Cowie, E. Douglas-Cowie, N. Tsapatsoulis, G. Votsis, S. Kollias, W. Fellenz, and J. Taylor, “Emotion recognition in human-computer interaction,” *IEEE Signal Processing Magazine*, vol. 18, no. 1, p. 3280, 2001.
- [9] Y. Tian, T. Kanade, and J. Cohn, “Facial expression analysis,” in *Handbook of face recognition*, S. Li and A. Jain, Eds. Springer, 2004.
- [10] M. Pantic, A. Pentland, A. Nijholt, and T. S. Huang, “Human computing and machine understanding of human behavior: A survey,” in *Artificial Intelligence for Human Computing*, ser. Lecture Notes in Artificial Intelligence, T. S. Huang, A. Nijholt, M. Pantic, and A. Pentland, Eds. London: Springer Verlag, 2007.
- [11] J. Wang, L. Yin, X. Wei, and Y. Sun, “3d facial expression recognition based on primitive surface feature distribution,” *IEEE Int’l Conf. on Computer Vision and Pattern Recognition*, vol. 2, pp. 1399–1406, 2006.
- [12] Y. Chang, M. Vieira, M. Turk, and L. Velho, “Automatic 3d facial expression analysis in videos,” *Analysis and Modelling of Faces and Gestures, Proceedings*, pp. 293–307, 2005.

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
- [13] B. Braathen, M. S. Bartlett, G. C. Littlewort, E. Smith, and J. R. Movellan, "An approach to automatic recognition of spontaneous facial actions," *Proc. Fifth IEEE Int'l Conf. Automatic Face and Gesture Recognition*, pp. 345–350, 2002.
  - [14] M. S. Bartlett, G. C. Littlewort, M. G. Frank, C. Lainscsek, I. R. Fasel, and J. R. Movellan, "Automatic recognition of facial actions in spontaneous expressions," *J. Multimedia*, vol. 1, no. 6, pp. 22–35, September 2006.
  - [15] F. Dornaika and F. Davoine, "Simultaneous facial action tracking and expression recognition using a particle filter," *Proc. Int'l Conf. on Computer Vision*, vol. 2, pp. 1733–1738, 2005.
  - [16] Y. Zhang and Q. Ji, "Active and dynamic information fusion for facial expression understanding from image sequences," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 5, pp. 699–714, May 2005.
  - [17] M. F. Valstar, I. Patras, and M. Pantic, "Facial action unit detection using probabilistic actively learned support vector machines on tracked facial point data," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition, Workshop Vision for Human-Computer Interaction*, June 2005.
  - [18] R. el Kaliouby and P. Robinson, "Real-time inference of complex mental states from facial expressions and head gestures," *Real-time vision for HCI*, pp. 181–200, 2005.
  - [19] B. Bascle and A. Blake, "Separability of pose and expression in facial tracking and animation," *Proc. Int'l Conf. on Computer Vision*, pp. 323–328, 1998.
  - [20] Z. Zhu and Q. Ji, "Robust real-time face pose and facial expression recovery," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 681–688, 2006.
  - [21] M. Anisetti, V. Bellandi, E. Damiani, and F. Beverina, "3d expressive face model-based tracking algorithm," *Proc. Signal Processing, Pattern Recognition, and Applications*, pp. 111–116, 2006.
  - [22] M. A. O. Vasilescu and D. Terzopoulos, "Multilinear analysis of image ensembles: Tensorfaces," *Proc. European Conf. on Computer Vision*, pp. 447–460, 2002.
  - [23] T. K. Marks, J. Hershey, J. C. Roddey, and J. R. Movellan, "Joint tracking of pose, expression, and texture using conditionally gaussian filters," *Advances in Neural Information Processing Systems*, vol. 17, pp. 889–896, 2005.
  - [24] A. Kapoor, Y. Qi, and R. W. Picard, "Fully automatic upper facial action recognition," *Proc. IEEE Int'l Workshop Analysis and Modeling of Faces and Gestures*, pp. 195–202, 2003.
  - [25] J. F. Cohn, L. I. Reed, Z. Ambadar, J. Xiao, and T. Moriyama, "Automatic analysis and recognition of brow actions and head motion in spontaneous facial behavior," *Proc. IEEE Int'l Conf. SMC*, vol. 1, pp. 610–616, 2004.
  - [26] J. F. Cohn and K. Schmidt, "The timing of facial motion in posed and spontaneous smiles," *Int'l J. Wavelets, Multiresolution and Information Processing*, vol. 2, pp. 1–12, March 2004.
  - [27] N. Sebe, M. Lew, I. Cohen, S. Yafei, T. Gevers, and T. Huang, "Authentic facial expression analysis," *Proc. Sixth IEEE Int'l Conf. Automatic Face and Gesture Recognition*, pp. 517–522, 2004.
  - [28] M. F. Valstar, M. Pantic, Z. Ambadar, and J. F. Cohn, "Spontaneous vs. posed facial behavior: Automatic analysis of brow actions," *Proc. Eighth Int'l conf. on Multimodal Interfaces*, pp. 162–170, 2006.
  - [29] G. Littlewort, M. S. Bartlett, and K. Lee, "Faces of pain: Automated measurement of spontaneous facial expressions of genuine and posed pain," *Proc. the Thirteenth Joint Symposium on Neural Computation*, p. 1, 2006.
  - [30] Z. Zeng, Y. Fu, G. Roisman, Z. Wen, Y. Hu, and T. S. Huang, "Spontaneous emotional facial expression detection," *J. of Multimedia*, vol. 1, no. 5, pp. 1–8, 2006.
  - [31] S. Ioannou, A. Raouzaoui, V. Tzouvaras, T. Mailis, K. Karpouzis, and S. Kollias, "Emotion recognition through facial expression analysis based on a neurofuzzy method," *Neural Networks*, vol. 18, no. 4, pp. 423–435, 2005.

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
- [32] S. Lucey, A. B. Ashraf, and J. Cohn, "Investigating spontaneous facial action recognition through aam representations of the face," in *Face Recognition Book*, K. Kurihara, Ed. Mammendorf, Germany: Pro Literatur Verlag, April 2007.
- [33] J. Russell and J. Fernandez-Dols, *The Psychology of Facial Expression*. New York: Cambridge Univ. Press, 1997.
- [34] J. N. Bassili, "Emotion recognition: The role of facial movement and the relative importance of upper and lower areas of the face," *J. Personality and Social Psychology*, vol. 37, no. 11, pp. 2049–2058, 1979.
- [35] in *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*, P. Ekman and E. Rosenberg, Eds. Oxford, UK: Oxford University Press, 2005.
- [36] Y. Tong, W. Liao, Z. Xue, and Q. Ji, "A unified probabilistic framework for facial activity modeling and understanding," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2007.
- [37] P. Ekman, W. V. Friesen, and J. C. Hager, *Facial Action Coding System: the Manual*. Salt Lake City, UT: Research Nexus, Div., Network Information Research Corp., 2002.
- [38] P. Ekman, *Telling Lies: Clues to Deceit in the Marketplace, Politics, and Marriage*. New York: Norton, 1985.
- [39] K. Schmidt and J. Cohn, "Dynamics of facial expression: Normative characteristics and individual differences," *IEEE Int'l Conf. on Multimedia and Expo*, pp. 728–731, 2001.
- [40] S. Nishio, K. Koyama, and T. Nakamura, "Temporal differences in eye and mouth movements classifying facial expressions of smiles," *Proc. Third IEEE Int'l Conf. Automatic Face and Gesture Recognition*, pp. 206–211, April 1998.
- [41] T. Dean and K. Kanazawa, "Probabilistic temporal reasoning," *AAAI88*, pp. 524–528, 1988.
- [42] T. Kanade, J. F. Cohn, and Y. Tian, "Comprehensive database for facial expression analysis," *Proc. Fourth IEEE Int'l Conf. AUTomatic Face and Gesture Recognition*, pp. 46–53, 2000.
- [43] "Intelligent systems lab (isl) database," <http://www.ecse.rpi.edu/homepages/cvrl/database/database.html>.
- [44] E. Douglas-Cowie, R. Cowie, and M. Schroeder, "The description of naturally occurring emotional speech," *Fifteenth Int'l Congress of Phonetic Sciences*, 2003.
- [45] M. A. of Discourse research lab, "<http://madresearchlab.org/>"
- [46] G. Schwarz, "Estimating the dimension of a model," *The Annals of Statistics*, vol. 6, pp. 461–464, 1978.
- [47] D. Heckerman, D. Geiger, and D. M. Chickering, "Learning bayesian networks: The combination of knowledge and statistical data," *Machine Learning*, vol. 20, no. 3, pp. 197–243, 1995.
- [48] S. J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*. Upper Saddle River, N.J.: Prentice Hall, 1995.
- [49] P. Wang and Q. Ji, "Multi-view face and eye detection using discriminant features," *Computer Vision and Image Understanding*, vol. 105, no. 2, pp. 99–111, February 2007.
- [50] Z. Zhu and Q. Ji, "Robust pose invariant facial feature detection and tracking in real-time," *Proc. Int'l Conf. Pattern Recognition*, vol. 1, pp. 1092–1095, 2006.
- [51] Y. Tong, Y. Wang, Z. Zhu, and Q. Ji, "Robust facial feature tracking under varying face pose and facial expression," *Pattern Recognition*, vol. 40, no. 11, pp. 3195–3208, November 2007.
- [52] K. Murphy, "Inference and learning in hybrid bayesian networks," *Technical Report CSD-98-990, Department of Computer Science, U. C. Berkeley*, 1998.
- [53] Y. Tong and Q. Ji, "Learning bayesian networks with qualitative constraints," *CVPR08*, 2008.
- [54] C. P. de Campos, Y. Tong, and Q. Ji, "Exploiting qualitative constraints for learning bayesian network parameters," *ECCV08*.
- [55] K. Murphy, "The bayes net toolbox for matlab," *Computing Science and Statistics*, vol. 33, 2001.