

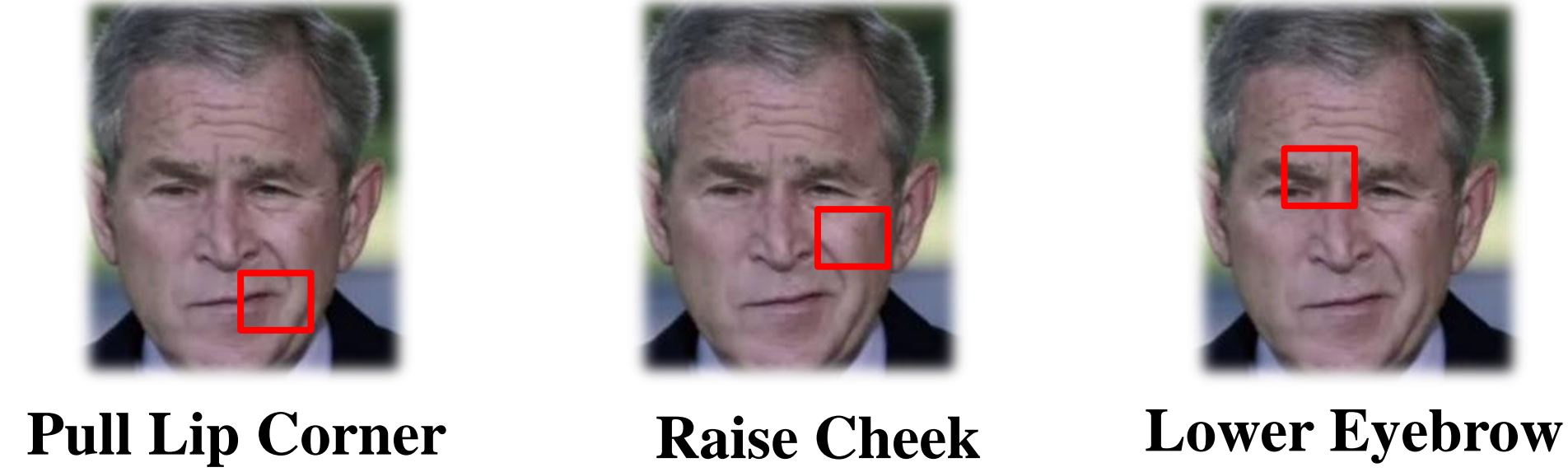
# Capturing Global Semantic Relationships for Facial Action Unit Recognition

Ziheng Wang<sup>1</sup> Yongqiang Li<sup>3</sup> Shangfei Wang<sup>2</sup> Qiang Ji<sup>1</sup>

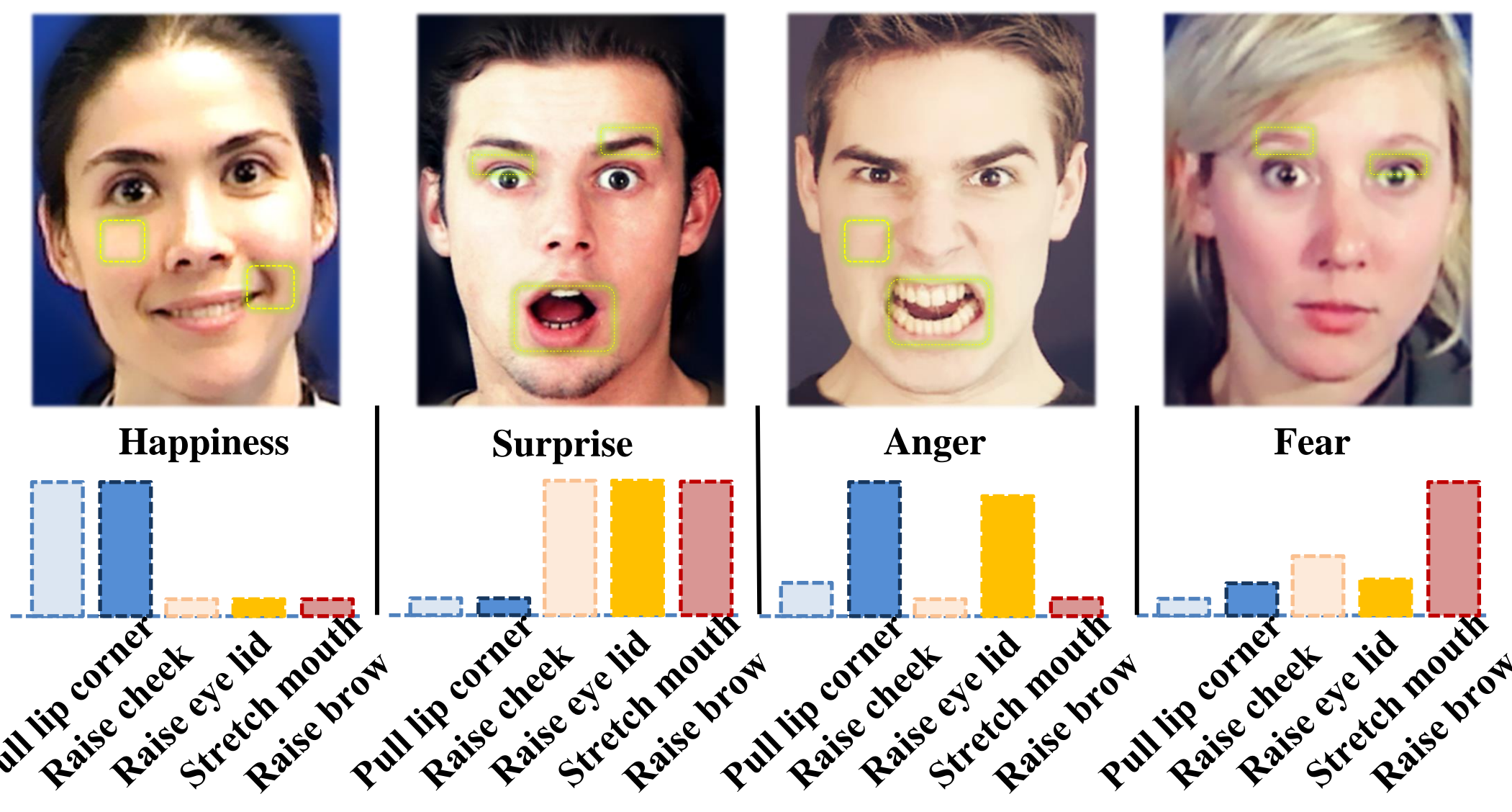
Rensselaer Polytechnic Institute<sup>1</sup> University of Science and Technology of China<sup>2</sup> Harbin Institute of Technology<sup>3</sup>

## 1. Problem

➤ Facial Action Unit Recognition



## 2. Main Idea

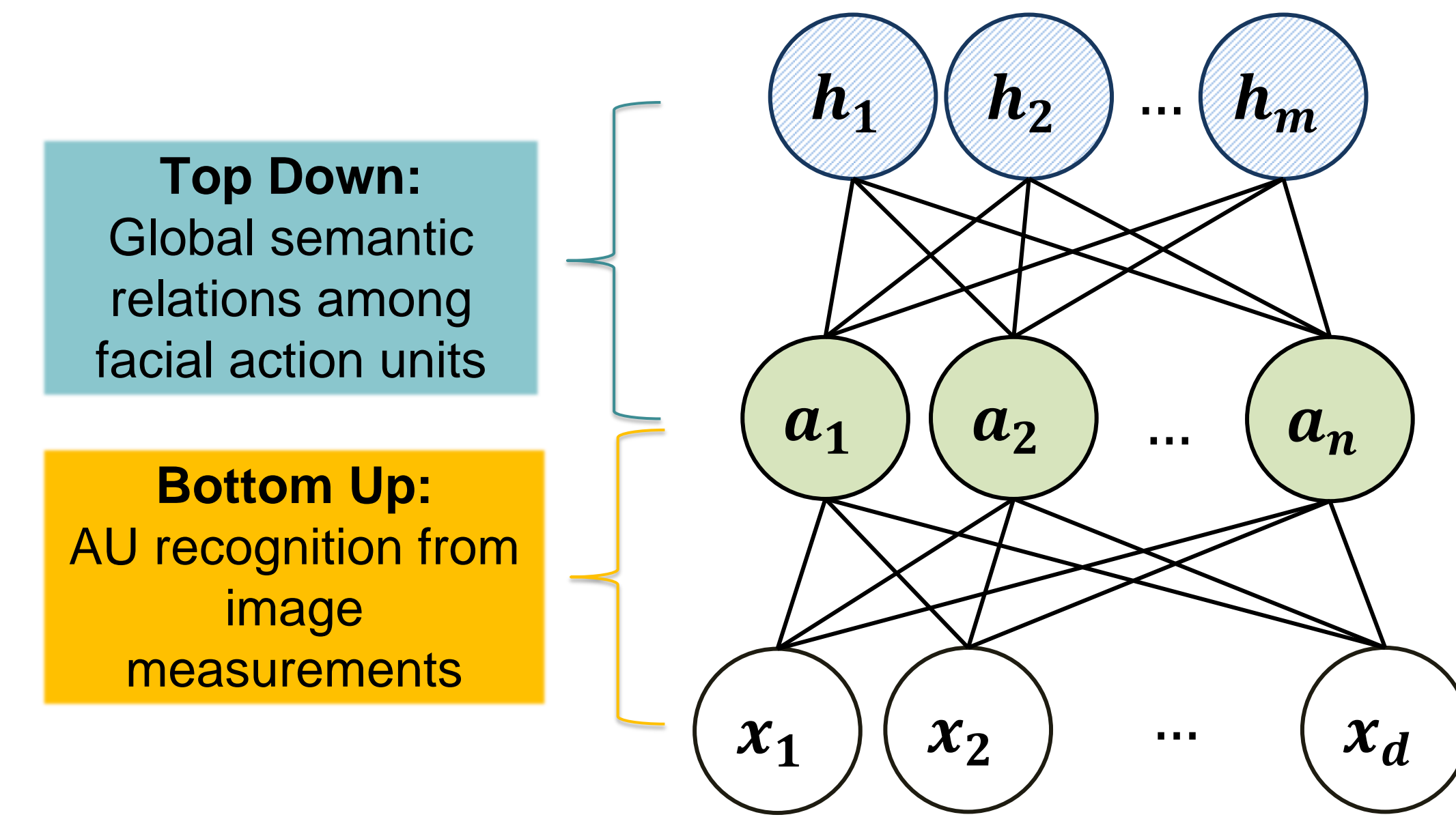


- Facial action units are NOT independent
  - AUs do not occur alone and some combinations of action units are frequently observed
  - Some AUs must or must not be present at the same time due to the limitations of facial anatomy
- Relationships among AUs are influenced by facial expression
  - “Stretch mouth” and “raise brow” are likely to be both absent during happiness, both present during surprise, and mutually exclusive during anger or fear
- Propose to use restricted Boltzmann machine (RBM) to capture global complex relationships among AUs for AU recognition
- Propose to use 3-way RBM to capture facial expression to more accurately characterize AU relationships

## 3. Related Work

- Existing approaches
  - Treat AUs are uncorrelated entities
  - Use Bayesian network (BN) to model AU relationships
- Limitations
  - Models such as BN are based on the first-order Markov assumption and therefore can only capture local, i.e. *pairwise* relationships between action units
  - Finding the optimal structure of a large AU network is difficult
  - Modeling AU relationships without considering the influence of facial expression, which could lead to incorrect estimation of AU dependencies
- Proposed Method
  - Model *global* AU relationships and consider the influence of *expression*

## 4. Hierarchical Model for AU Recognition

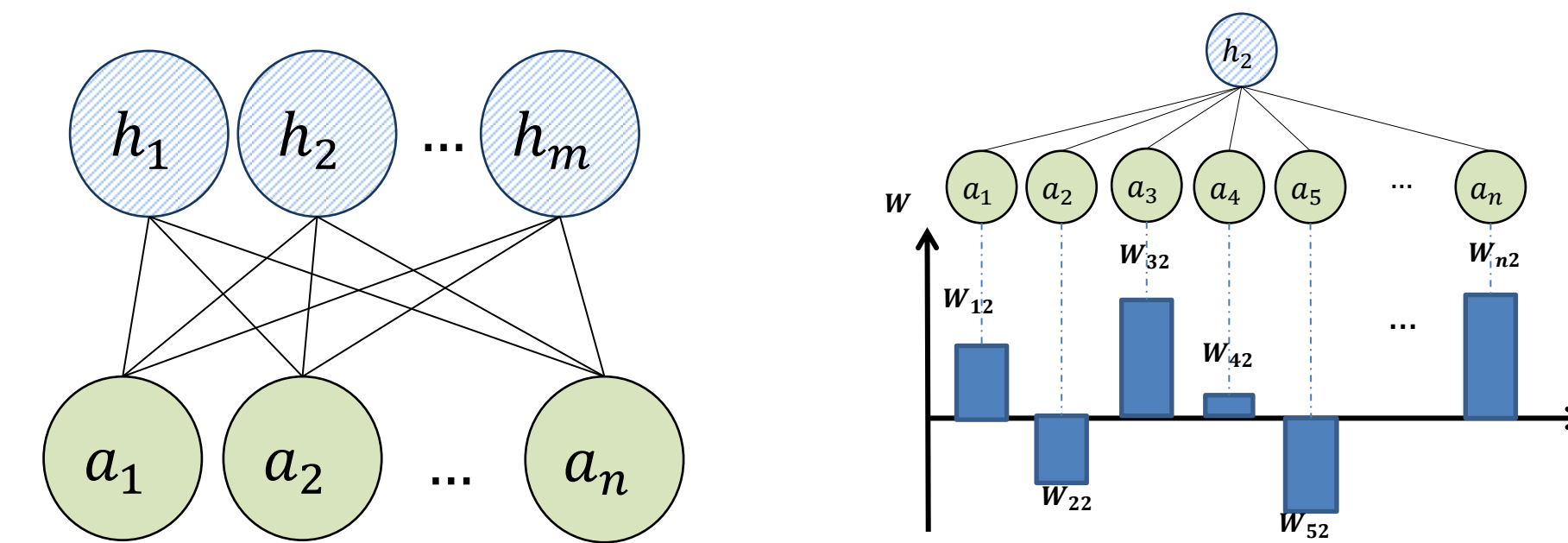


- Middle layer  $a_1$  to  $a_n$ : binary state of  $AU_1$  to  $AU_n$
- Bottom layer  $x_1$  to  $x_d$ : image features
- Top layer  $h_1$  to  $h_m$ : latent nodes modeling AU relationships
- Total Energy:

$$E(x, a, h; \theta) = -\sum_i \sum_j a_i W_{ij}^1 h_j - \sum_j c_j h_j - \sum_i b_i a_i - \sum_i \sum_t W_{it}^2 a_i x_t$$

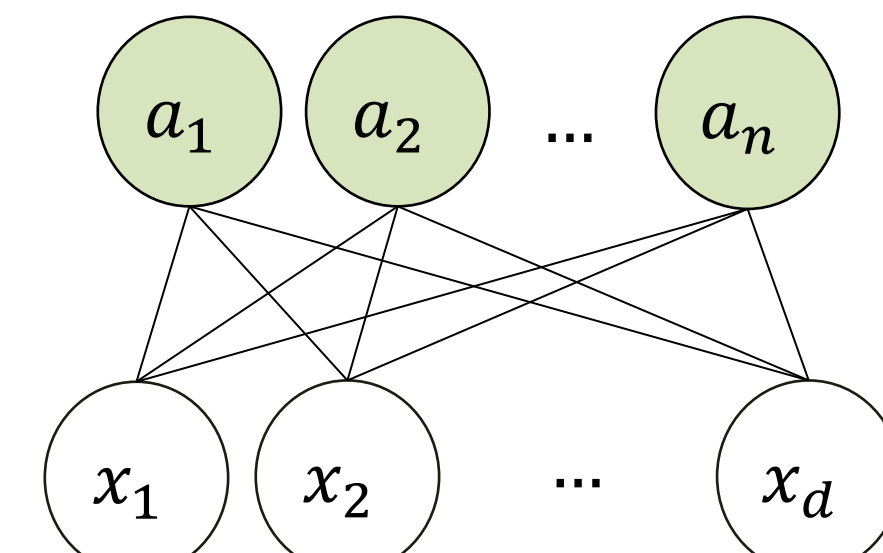
- $a_i W_{ij}^1 h_j$ : compatibility between  $AU_i$  and latent node  $h_j$
- $c_j h_j$ : bias for each latent node  $h_j$
- $b_i a_i$ : bias for each AU node  $a_i$
- $W_{it}^2 a_i x_t$ : compatibility between  $AU_i$  and feature  $x_t$

➤ Top down - capturing global relations among AUs



- Each latent node is connected to all AU nodes and therefore modeling their higher-order relationships
- The captured AU relationships can be implicitly inferred from the model parameters  $W_{ij}^1$
- Vector  $[w_{im}]_{i=1}^n$  captures a specific presence and absence pattern of all the action units
- $W_{im}^1$  large  $\rightarrow AU_i$  more likely to occur in pattern  $m$
- $W_{im}^1$  small  $\rightarrow AU_i$  less likely to occur in pattern  $m$

➤ Bottom up: AU recognition from image features



- Each AU node is connected to all the image features with the energy  $E(a_i, x) = -\sum_t W_{it}^2 a_i x_t$
- Equivalent to a set of linear AU classification models

## 5. Learning and Inference

➤ Discriminative Learning

Given training data  $\{(x_i, a_i)\}_{i=1}^N$

$$\theta^* = \arg \max_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log P(a_i | x_i; \theta)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{P(h|a, x; \theta)} - \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{P(h, a | x; \theta)}$$

- Calculating the gradient requires  $P(h|a, x; \theta)$  and  $P(h, a | x; \theta)$
- $P(h|a, x; \theta)$  can be analytically computed

$$P(h_j | a, x; \theta) = P(h_j | a; \theta) = \sigma \left( -c_j - \sum_i W_{ij}^1 a_i \right) \quad (1)$$

- $P(h, a | x; \theta)$  is intractable, we revised contrastive divergence (CD) algorithm to compute it
- The basic idea is to approximate  $P(h, a | x; \theta)$  by sampling  $h$  with Equation 1 and then sampling  $a$  with Equation 2

$$P(a_i | h, x; \theta) = \sigma \left( -b_i - \sum_j W_{ij}^1 h_j - \sum_t W_{it}^2 x_t \right) \quad (2)$$

➤ Inference

- Given a query sample  $x$ , each action unit can be inferred by

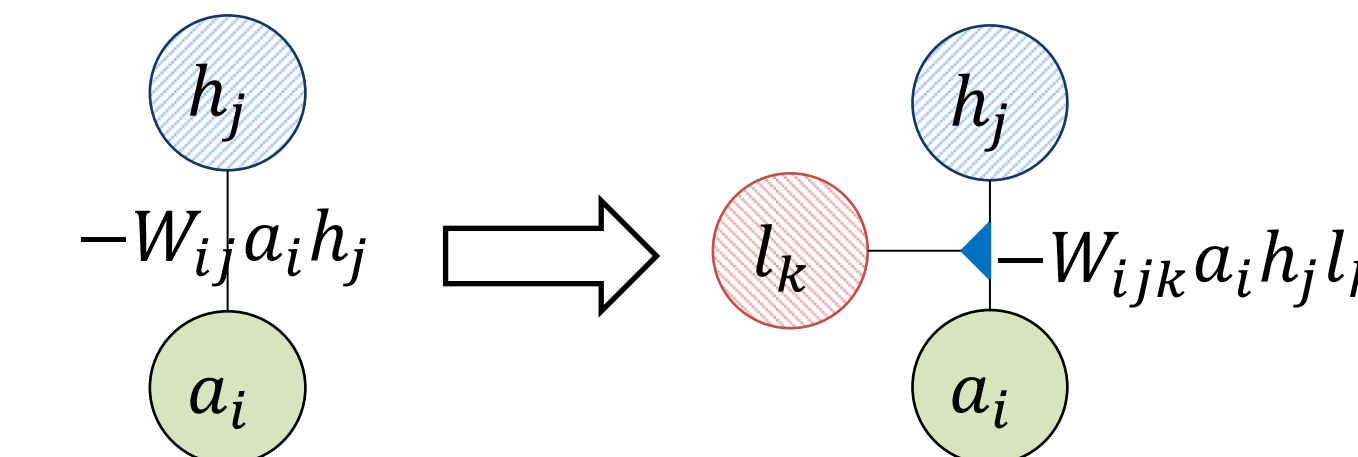
$$a_i^* = \arg \max_{a_i} P(a_i | x)$$

- Can be efficiently performed with Gibbs sampling by iteratively sampling  $h$  from  $P(h|x, a)$  and sampling  $a$  from  $P(a|h, x)$

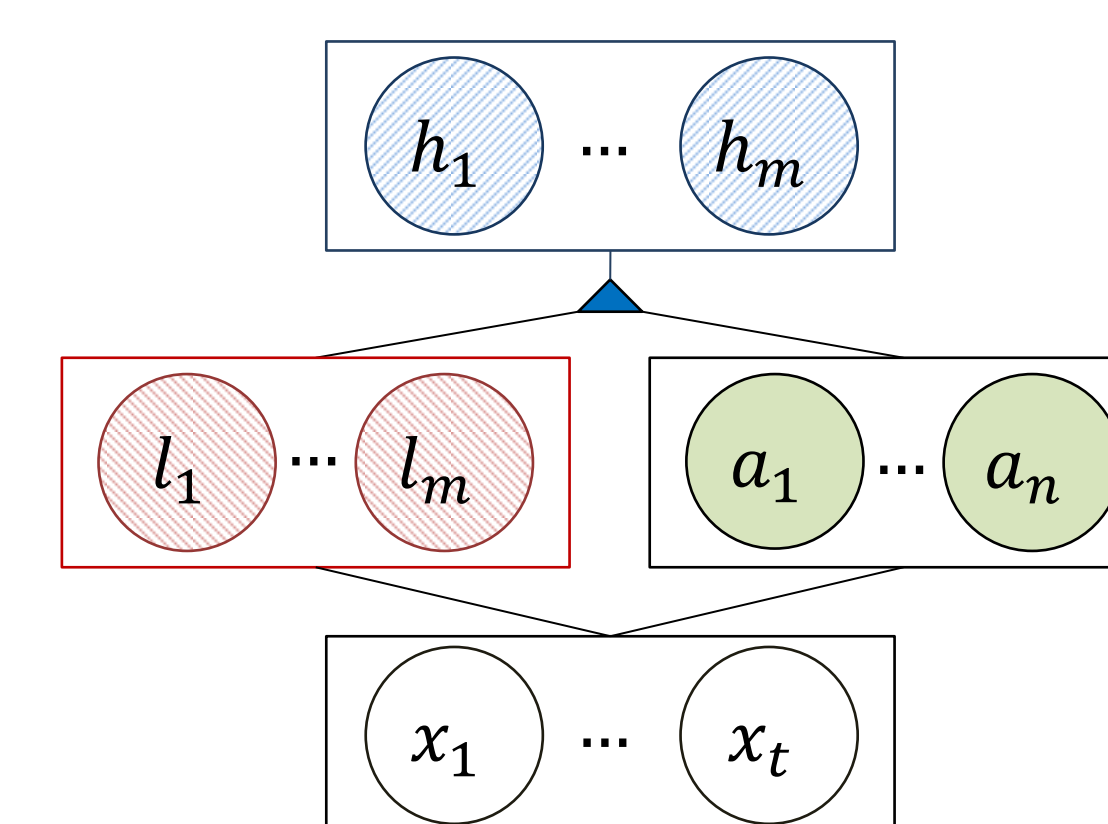
## 6. Incorporating Expression to Model AU Relations

➤ Relations among the AUs depend on facial expressions

- Expression is known during training but known during testing
- Basic idea: modulate the connection between each pair of action unit and latent unit ( $a_i, h_j$ ) with an expression variable  $l$



○ Model



$$E(x, a, l; \theta) = -\sum_i \sum_j \sum_k W_{ijk}^1 a_i h_j l_k - \sum_i \sum_t W_{it}^2 a_i x_t - \sum_k \sum_t W_{kt}^3 l_k x_t - \sum_j c_j h_j - \sum_i b_i a_i - \sum_k d_k l_k$$

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log P(a_i, l_i | x_i; \theta)$$

$$a_i^* = \arg \max_{a_i} P(a_i | x) = \arg \max_{a_i} \sum_l P(a_i, l | x)$$

## 7. Experimental Results

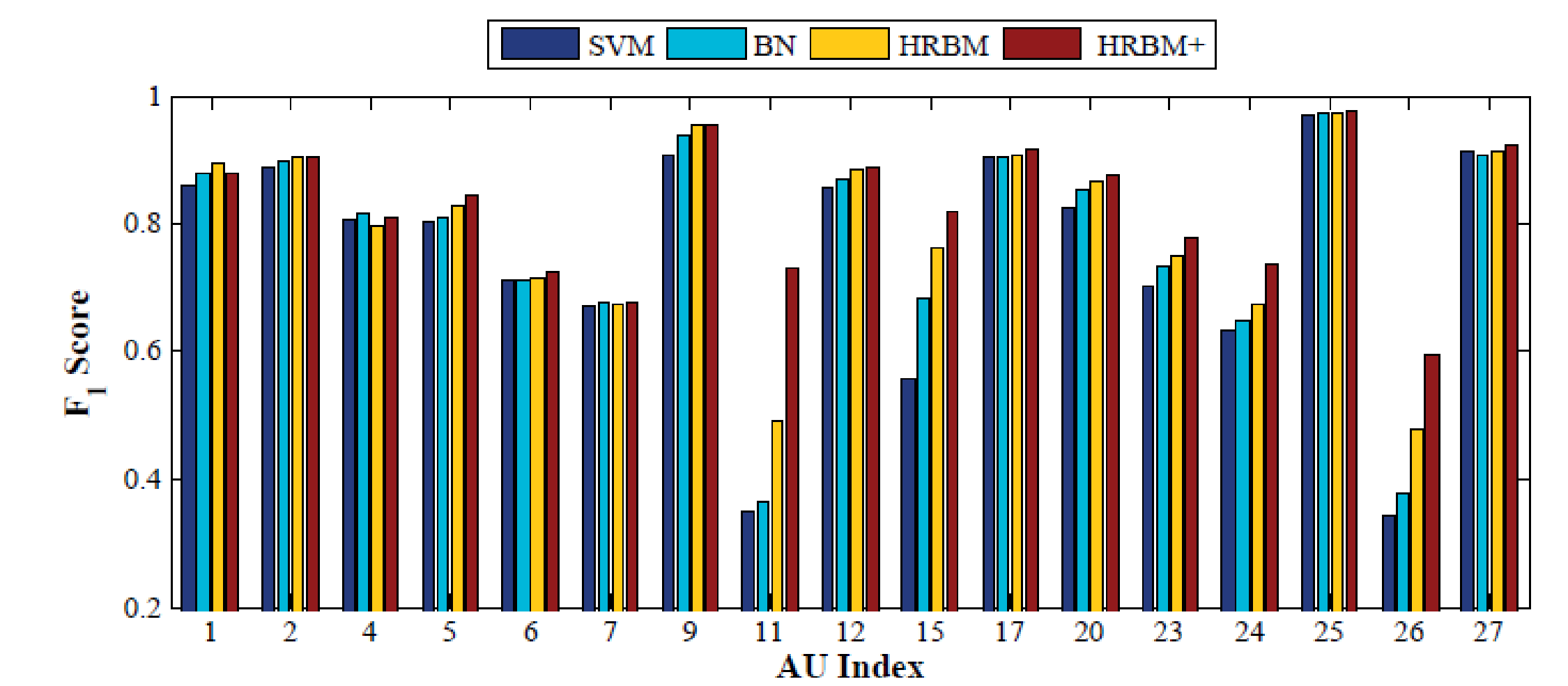
➤ Methods

- SVM, BN [Tong *et al.* 2007], HRBM, HRBM+

➤ Posed Facial Action Units Recognition

- CK+: 593 peak images, 17 action units

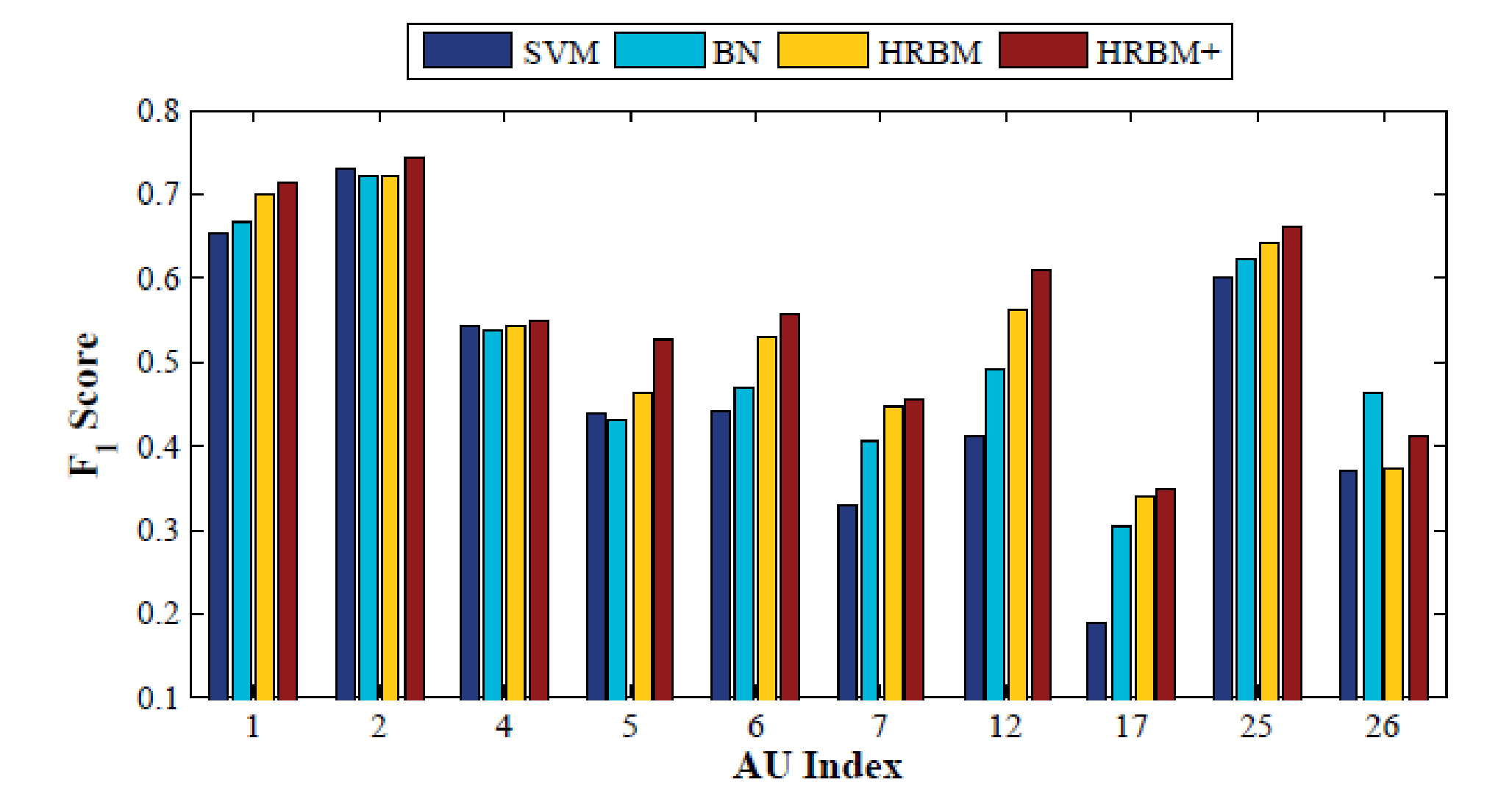
Method	SVM	BN	HRBM	HRBM+
Average $F_1$ -score	74.70%	76.70%	79.21%	82.44%



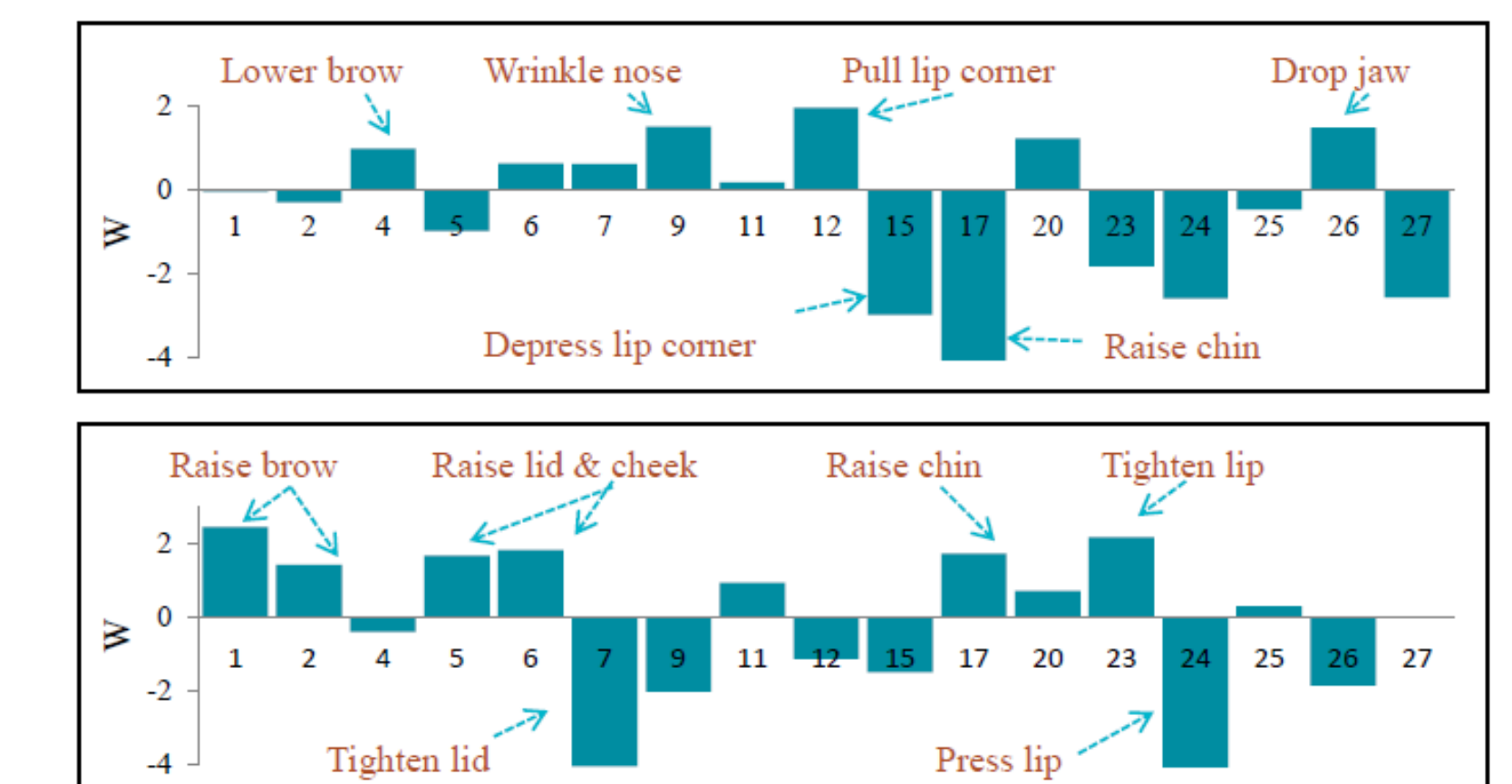
➤ Non-posed Facial Action Unit Recognition

- SEMAIN: 180 image frames, 10 action units
- Train on CK+, test on SEMAIN

Method	SVM	BN	HRBM	HRBM+
Average $F_1$ -score	47.70%	51.09%	54.76%	56.14%



➤ Semantic Relationship Analysis



## 8. Conclusions

- Proposed a hierarchical model for AU recognition
- Capture higher-order AU interactions
- Consider the influence of facial expression on AU relations
- Experimental results demonstrate the effectiveness of the proposed approach