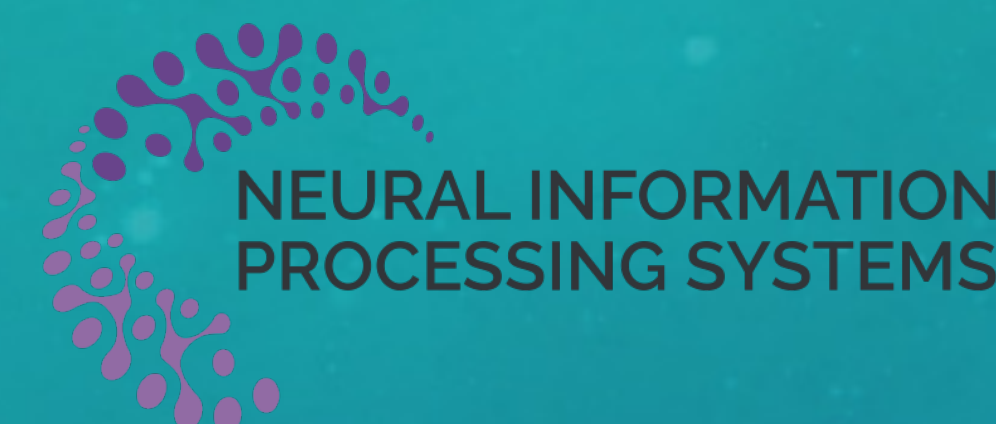# Knowledge Augmented Deep Neural Networks for Joint Facial Expression and Action Unit Recognition

Zijun Cui (`cuiz3@rpi.edu`)
Rensselaer Polytechnic Institute

Tengfei Song
Southeast University

Yuru Wang
Northeast Normal University

Qiang Ji (`jiq@rpi.edu`)
Rensselaer Polytechnic Institute

**NEURAL INFORMATION PROCESSING SYSTEMS**

## Introduction



Expression: 'HAPPY'
AU6(Cheek Raiser): 'ON'
AU12(Lip Corner Puller): 'ON'

Figure: An example from CK+ dataset[1]

❑ Tasks:
- Facial Expression Recognition(FER)
- Action Unit(AU) Detection

❑ Motivations:
- Facial expression and AUs are strongly correlated
- Generic knowledge on expression-AUs relationships is available

❑ Contributions:
- A *knowledge model* encoding the generic knowledge systematically
- A deep learning framework for *joint* facial expression and AU recognition

## Generic Knowledge as Probabilities
-- on expression-AUs probabilistic relationships

❑ Notation:
- Expression $X^e = \{1,2,...,E\}$
  E is the total number of expressions
- AUs $X^{au}_m = \{X^{au}_1, X^{au}_2, ..., X^{au}_M\}$
  M is the total number of AUs and $X^{au}_m = \{0,1\}$

❑ *Expression-dependent single AU probabilities*
  o AU4 is a primary AU given Anger expression
$$p(X^{au}_4 = 1|X^e = Anger) > p(X^{au}_4 = 0|X^e = Anger)$$

❑ *Expression-dependent joint AU probabilities*
  o AU6 and AU12 are positively correlated given Happy expression
$$p(X^{au}_6 = 1|X^{au}_{12} = 1, X^e = Happy) > p(X^{au}_6 = 0|X^{au}_{12} = 1, X^e = Happy)$$
$$p(X^{au}_6 = 1|X^{au}_{12} = 1, X^e = Happy) > p(X^{au}_6 = 1|X^{au}_{12} = 0, X^e = Happy)$$

❑ *Expression-independent joint AU probabilities*
  o AU1 and AU2 are positively correlated
$$p(X^{au}_1 = 1|X^{au}_2 = 1) > p(X^{au}_1 = 0|X^{au}_2 = 1)$$
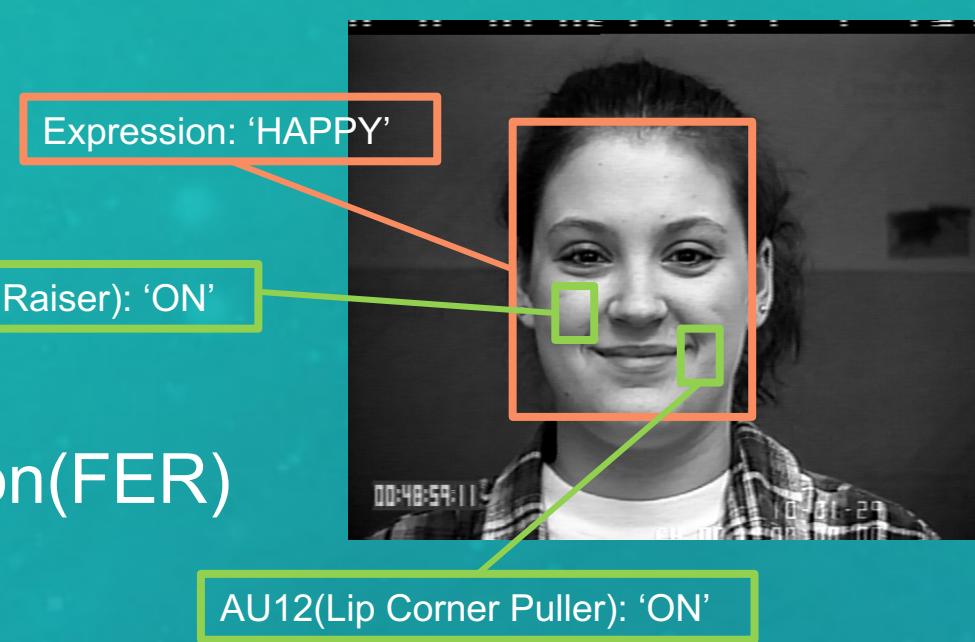$$p(X^{au}_1 = 1|X^{au}_2 = 1) > p(X^{au}_1 = 1|X^{au}_2 = 0)$$

## Encoding of the Generic Knowledge
-- Bayesian Network(BN) Learning with Probability Constraints



Figure: Overview of the proposed framework

❑ Definitions of the Bayesian Network
- Conditional probabilities are parameterized with the regression equations
$$p(X_i = k|\pi(X_i)) = \sigma_M(\sum_{j=1}^J w_{ijk}\pi_j(X_i) + b_{ik})$$
where weights $w = \{w_{ijk}\}$ and bias $b = \{b_{ik}\}$ are to be learned
- $A(w)$ is the weighted adjacency matrix defining the structure[2]: $A_{ij} = \sum_{k=1}^K ||w_{ijk}||^2_2$
- The constraint of Directed Acyclic Graph(DAG)[3]: $\text{tr}(e^{A(w)\circ A(w)}) - N = 0$

❑ Probability constraints derived from the generic knowledge
- Strictly inequality constraints: $\{g_i(w,b) < 0\}_{i=1}^G$
  To better handle $g_i$, we define positive margin with additional variable $s_i$
  The strictly inequality constraints become equality constraints:
$$g_i(w,b) + e^{s_i} = 0, i = 1,...,G$$
- Inequality constraints: $\{l_j(w,b) \le 0\}_{j=1}^L$
- Equality constraints: $\{h_k(w,b) = 0\}_{k=1}^H$
- For example:
$$g_i(w,b) = p(X^{au}_1 = 0|X^{au}_2 = 1; w,b) - p(X^{au}_1 = 1|X^{au}_2 = 1; w,b) < 0$$

❑ A penalty function $f(w,b;s)$ measures the violation of constraint
$$f(w,b;s) = \frac{1}{G}\sum_{i=1}^G \log((g_i(w,b) + e^{s_i})^2 + 1)$$
$$+ \frac{1}{L}\sum_{j=1}^L \log\left(\left(l_j^+(w,b)\right)^2 + 1\right) + \frac{1}{K}\sum_{k=1}^K \log\left(\left(h_i(w,b)\right)^2 + 1\right)$$
with weights $w$, bias $b$ and current margins $e^s$. And $l_j^+ = \max\{0, l_j\}$
- $f(w,b;s) = 0$ if and only if all the constraints are satisfied

❑ A *Constrained Optimization Approach* for BN learning
$$w^*, b^*, s^* = \arg\min_{w,b,s} f(w,b;s) + \gamma||w||_1 - \mu||s||^2_2$$
$$s.t. \ \text{tr}(e^{A(w)\circ A(w)}) - N = 0$$
where $||w||_1$ penalizes the density of the structure, and $||s||^2_2$ encourages the bigger positive margins

❑ The learned Bayesian Network serves as our knowledge model $K$

## AU detection model and FER models

❑ We learn AU detection model and FER models with:
- The training images $x_n, n = 1, ..., N$
- The GT expression labels $y^{GT}_n, n = 1, ..., N$
- The knowledge model $K$
* $N$ is the total number of training samples

❑ Phase 1: Initialization of AU detection and FER models
  ➤ Weakly supervised AU detection model $g_\varphi$
$$\varphi^* = \text{argmin}_\varphi \frac{1}{N}\sum_{n=1}^N E_{p(z_n|y^{GT}_n, K)} l(z_n, g_\varphi(x_n))$$
  $p(z_n|y^{GT}_n, K)$ is the probability of AU configuration $z_n$ computed from the BN model and the $y^{GT}_n$
  ➤ Facial Expression Recognition(FER) Models
- Image-based FER model $f_\psi$: $\psi^* = \text{argmin}_\psi \frac{1}{N}\sum_{n=1}^N l(y^{GT}_n, f_\psi(x_n))$
- AU-based FER model $h_\phi$: $\phi^* = \text{argmin}_\phi \frac{1}{N}\sum_{n=1}^N l(y^{GT}_n, h_\phi(g_\varphi(x_n)))$
  where $g_\varphi(x_n)$ is the output of the AU model $g_\varphi$
* $l$ is the cross-entropy loss

❑ Phase 2: Integration among AU and Expression Models
  ➤ The combined expression probability
$$p(y_n|x_n, K) = w_1 p_\psi(y_n|x_n) + w_2 p_\phi(y_n|g_\varphi(x_n), K)$$
  $p_\psi(y_n|x_n)$ is the output of $f_\psi$ and $p_\phi$ is the output of $h_\phi$. $w_1, w_2$ are the weights

  ➤ Expression-augmented AU detection model
$$\varphi^* = \text{argmin}_\varphi \frac{1}{N}\sum_{n=1}^N E_{p(z_n|y^{GT}_n, K)} l\left(z_n, g_\varphi(x_n)\right) + \lambda_1 E_{p(y_n|x_n, K)}E_{p(z_n|y_n, K)} l\left(z_n, g_\varphi(x_n)\right)$$

  ➤ Knowledge-augmented image-based FER model
$$\psi^* = \text{argmin}_\psi \frac{1}{N}\sum_{n=1}^N l(y^{GT}_n, f_\psi(x_n)) + \lambda_2 E_{p(y_n|x_n, K)} l\left(y_n, f_\psi(x_n)\right)$$

* $l$ is the cross-entropy loss. $\lambda_1, \lambda_2$ are the hyper-parameters to be tuned

## Experiments
-- comparisons with state-of-the-art models

- Action Unit Detection

Table 6: Comparison to the SoAs on AU detection.

| Supervision | Method | BP4D | CK+ | MMI |
|---|---|---|---|---|
| Supervised | HRBM[47] | .67 | .79 | .56 |
| | MC-LVM[8] | - | .80* | - |
| | JPML[56] | .68* | .78* | - |
| | AU R-CNN[30] | .63* | - | - |
| Weakly-supervised | HTL[40] | .50 | .66 | .42 |
| | LP-SM[54] | .55 | .72* | .50 |
| | TCAE[22] | .56* | - | - |
| | AUD-BN(baseline) | .56 | .69 | .47 |
| | AUD-EA(gBN) | .57 | .74 | .58 |

- Facial Expression Recognition(FER)

Table 8: Comparison with SoA FER methods

| Methods | BP4D | CK+ | MMI | EmotioNet |
|---|---|---|---|---|
| STM-Explet[27] | - | 94.19* | 75.12* | - |
| DTAGN(Joint)[12] | - | 97.25* | 70.24* | - |
| DeRL[50] | - | 97.30* | 73.23* | - |
| ILCNN[3] | - | 94.35* | 70.67* | - |
| DAM-CNN[49] | - | 95.88* | - | - |
| FMPN-FER[4] | 60.16 | 96.53 | 82.74* | 84.88 |
| DeepEmotion[32] | 79.54 | 95.23 | 72.66 | 81.51 |
| FER-I(baseline) | 61.68 | 94.29 | 67.35 | 80.85 |
| FER-IK(gBN) | 83.82 | 97.59 | 84.90 | 95.55 |

[1] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews. The extended cohn-kanadedataset (ck+): A complete dataset for action unit and emotion-specified expression. CVPR 2010
[2] Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, and Eric Xing Learning sparse nonparametricdags. IInternational Conference on Artificial Intelligence and Statistics, 2020
[3] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing Dags with no tears: Continuousoptimization for structure learning. InAdvances in Neural Information Processing Systems, 2018.