

AU-Expression Knowledge Constrained Representation Learning for Facial Expression Recognition

Tao Pu¹, Tianshui Chen^{2,*}, Yuan Xie², Hefeng Wu¹, and Liang Lin^{1,2}

Abstract—Recognizing human emotion/expressions automatically is quite an expected ability for intelligent robotics, as it can promote better communication and cooperation with humans. Current deep-learning-based algorithms may achieve impressive performance in some lab-controlled environments, but they always fail to recognize the expressions accurately for the uncontrolled in-the-wild situation. Fortunately, facial action units (AU) describe subtle facial behaviors, and they can help distinguish uncertain and ambiguous expressions. In this work, we explore the correlations among the action units and facial expressions, and devise an AU-Expression Knowledge Constrained Representation Learning (AUE-CRL) framework to learn the AU representations without AU annotations and adaptively use representations to facilitate facial expression recognition. Specifically, it leverages AU-expression correlations to guide the learning of the AU classifiers, and thus it can obtain AU representations without incurring any AU annotations. Then, it introduces a knowledge-guided attention mechanism that mines useful AU representations under the constraint of AU-expression correlations. In this way, the framework can capture local discriminative and complementary features to enhance facial representation for facial expression recognition. We conduct experiments on the challenging uncontrolled datasets to demonstrate the superiority of the proposed framework over current state-of-the-art methods. Codes and trained models are available at <https://github.com/HCP Lab-SYSU/AUE-CRL>.

I. INTRODUCTION

Facial expression recognition (FER) is essential for intelligent robotics because it can help the robotics to understand human emotions and behaviors. Basically, this task aims to classify basic (e.g., happy, angry) or compound (e.g., happy & surprised, sad & angry) expressions based on face appearance for both in-the-lab [1], [2], [3] and in-the-wild environments [4], [5]. Recently, most fruitful algorithms resort to deep neural networks [6], [7] to learning powerful feature representation to promote the FER performance. Despite achieving impressive progress for the lab-controlled environments, it is still challenging and unsolved due to the complex variations in pose, illumination, and age, especially in the uncontrolled environments during the natural human robot interaction process.

According to the facial action coding system (FACS) [8], [9], facial action units (AUs) encode subtle facial appearances and changes, which have strong correlations

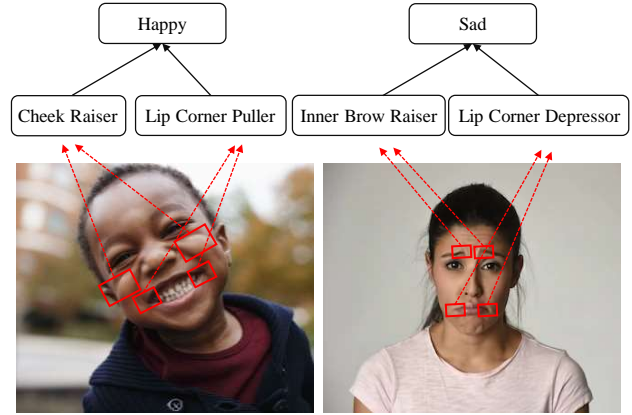


Fig. 1. Two examples of the correlations between facial expression and action units.

with human expressions. For example as shown in Figure 1, if a face image is detected to have the AUs of “cheek raiser” and “lip corner puller”, it tends to be a “happy” face. In contrast, if the detected AUs are “inner brow raiser” and “lip corner depressor”, it is more likely to be a “sad” face. Thus, automatically detecting the AUs and modeling their relationships with the expressions is essential to promote FER performance, especially to help distinguish uncertain or ambiguous expressions.

In this work, we aim to mine the AU features to enhance face image representation for more robust and accurate expression recognition. To achieve this end, two crucial challenges arise. First, most current FER datasets (e.g., RAF-DB [4], SFEW2.0 [5]) do not have AU annotations, and it is very expensive and labor-consuming to annotate the AUs for these datasets. Thus, how to learn AU features without AU annotations is a key challenge. Second, there exist tens of AUs, and not all AUs are equally important for different expressions. How to adaptively select AU features to enhance image representation for each expression is another vital problem.

To address these challenges, we explore exploiting the correlations among expressions and AUs to learn AU features in an unsupervised manner and adaptively select these features for feature enhancement by developing a novel AU-Expression Knowledge Constrained Representation Learning (AUE-CRL) framework. Specifically, there exist strong correlations among expressions and AUs. We first design a knowledge-guided AU representation learning module that leverages these correlations to covert expression labels to pseudo AU labels and utilizes pseudo labels to train the AU classifiers to obtain

*Tianshui Chen is the corresponding author

¹Tao Pu, Tianshui Chen, Hefeng Wu, and Liang Lin are with Human Cyber Physical Intelligence Integration Lab, School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China {putao537, tianshuichen, wuhefeng}@gmail.com

²Tianshui Chen, Yuan Xie, and Liang Lin are with DarkMatter AI Research, Guangzhou, China. {tianshuichen, phoenixsysu}@gmail.com

the AU features. As suggested in previous work, different AUs also have obvious co-occurrence dependencies, and these dependencies are important for selecting useful AUs. With the AU features, we further introduce a knowledge-guided attention module that learns to adaptively mine the most relevant AU features under the constraint of AU dependencies. In this way, the framework can automatically discover useful local facial behaviors to facilitate FER.

In summary, the contributions of this work are three-fold. First, we design a novel AU-Expression Knowledge Constrained Representation Learning (AUE-CRL) framework that exploits prior knowledge of AU-expression correlations to automatically discover useful AU features for feature enhancement to facilitate recognizing facial expression. Second, we propose to leverage the AU-expression correlations to guide learning AU features without AU annotations. Thus, the framework can be easily generalized to all of current FER datasets. Finally, we conduct extensive experiments on several large-scale in-the-wild datasets to demonstrate the effectiveness of the proposed framework, and carry out ablative studies to analyze the actual contribution of each key component.

II. RELATED WORK

In this section, we review the most related works about facial expression recognition and facial action unit detection.

A. Facial Expression Recognition

Previous works on facial expression recognition mainly focused on the basic categories (e.g., happy, angry, etc.) in which the data were collected by asking volunteers to make specific expressions in the constrained lab environments [10], [11]. Traditional methods primarily designed hand-crafted features (e.g., Local Binary Pattern (LBP) [10], Bag of Word (BoW) [12], and Histogram of Oriented Gradient (HoG) [11]). These methods can achieve satisfactory performance in such constrained environments. Recently, there emerged various large-scale datasets, in which the data were captured in real-world scenarios [13], [14], [15]. Compared with previous in-the-lab settings, these datasets were even more challenging due to more variance in pose, illumination, etc. Previous hand-crafted features could hardly capture such variance and thus they worked quite poor on these datasets. To address this issue, recent works resorted to deep convolutional networks [6], [7] to learn more powerful feature representation for expression recognition. For example, Liu et al. [15] proposed a Boosted Deep Belief Network to learn feature that could characterize expression-related facial appearance/shape changes. Mollahosseini et al. [16] designed deeper convolutional networks that built on inception module [17] to extract more discriminative feature for recognition. Liu et al. [18] exploited 3D Convolutional Neural Networks that was trained with deformable action parts constraints to learn discriminative part-based representation. Yu et al. [19] ensembled multiple deep Convolutional Neural Networks to further improve recognition accuracy. Different from these

works, our method introduces the relationships between AU-expression and relationship among AUs to mine information of AUs to help capture local facial behaviors to facilitate facial expression recognition.

The proposed method is also related to some recent works that also exploit prior knowledge to facilitate visual reasoning [20], [21], [22], [23], [24], [25]. For example, Chen et al. introduce the co-occurrence correlations among different categories to help better recognize multiple semantic objects [23], [22]. Generally, these methods represent prior knowledge in the form of graphs and introduce the graph neural networks [26], [27] for message propagation to learn contextualized feature representations. Different from these works, we introduce the relationships between AU-expression and among AUs to guide generating pseudos AU labels for AU representation learning and selecting useful AUs for facial expression recognition.

B. Action Unit Recognition

Facial action units are defined to describe facial muscle movements according to [28], and detecting action unit is helpful for expression and emotion understanding. Previous works [29], [30] leveraged traditional shadow models (e.g., SVM and SVR) to solve this task. For example, Mahoor et al. [30] projected high-dimensional facial images into a low-dimensional space via the spectral regression technique and adopted SVM classifier to predict the AU intensity. Inspired by recent advance of deep neural network on vision tasks, current works [31], [32], [33], [34] also designed deep model for action unit recognition. Gudi et al. [34] designed a simple seven-layer network to estimate both occurrences and intensities of the AUs. Furthermore, Walecki et al. [33] adopted the conditional random field (CRF) to encode AU dependencies and combined it with deep neural networks to improve recognition.

III. AUE-CRL FRAMEWORK

A. Overview

This proposed AUE-CRL framework explores mining useful local AU information to enhance image representation learning under the guidance of the correlations among AUs and expressions. It mainly consists of modules, i.e., feature extractors, knowledge-guided AU representation learning (KGAURL) module, and knowledge-constrained AU selection (KCAUS) module. Taking an input image I , the feature extractor generates multi-layer feature maps, and then it fuses these feature maps and performs global average pooling to obtain global expression feature vector $\mathbf{f}^e \in \mathcal{R}^{d_e}$. The feature extractor also fuses the feature maps from multiple layers inversely to obtain a set of feature maps with relative large size and leverage a crop network to extract feature for each AU, i.e., $\{\mathbf{f}_1^a, \mathbf{f}_2^a, \dots, \mathbf{f}_A^a\}$. Here, $\mathbf{f}_i^a \in \mathcal{R}^{d_a}$ is the feature vector of the i -th AU and A is the AU number. Then, the KGAURL module converts the expression labels to pseudo AU labels based on the AU-expression correlations. The pseudo labels are then used to supervised AU classifier training and thus it can learn AU features without any additional AU annotations.

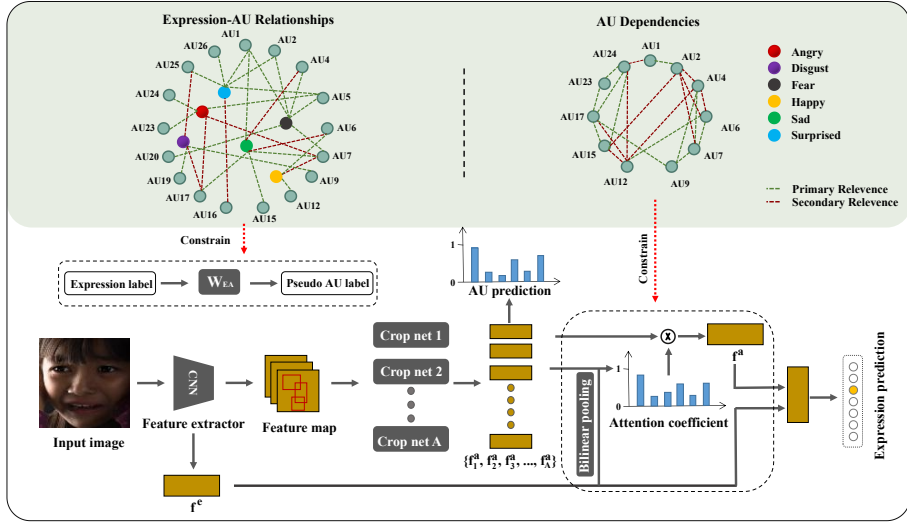


Fig. 2. An illustration of the proposed framework. We exploit the relationships between AUs and expressions among AUs to constrain the attention mechanism to better select useful AU representation to facilitate facial expression recognition. To keep this illustration concise, we only show part of relationships and dependencies.

Finally, the KCAUS exploits an attention mechanism to automatically discover useful AU features under the constraint of AU dependencies, and aggregate these features with the global expression feature vector for expression prediction. An overall framework is illustrated in Figure 2.

B. Knowledge-Guided AU Representation Learning

Current algorithms [13], [34] mainly resort to deep neural networks [7], [6], [35] to learn AU representations, but they require a large amount of ground truth annotations to ensure the discriminative and generalization abilities. However, current datasets with AU annotations are mainly captured in the constrained in-the-lab environment and cover very few subjects, e.g., 27 and 41 subjects on the DISFA [36] and BP4D [37] datasets. Features trained using these datasets can hardly generalize to other environments and subjects, especially to the in-the-wild settings. On the other hand, current FER datasets lack the AU annotations, and it is expensive and unpractical to add the AU annotations for these datasets. In this work, we design a knowledge-guided module that exploits the correlations between AU and expression to generate pseudo AU labels, and thus the proposed framework can learn AU representation without incurring additional annotations.

As suggested in previous literatures [28], each expression is relevant with several AUs, and the relevance can be further divided into primary and secondary ones. More concretely, if a face image is annotated with a specific expression, it tends to have the corresponding primary AUs with high probabilities, secondary AUs with middle probabilities, and other AUs with low probabilities. According to these relationships, we can build a correlation matrix $\widehat{\mathbf{W}}_{EA}$, where E and A denote the expression and AU number, respectively. The value \hat{w}_{ea} denotes the prior relevant probability between expression e and AU a . It is assigned with a large value if they are primarily

relevant, a middle value if they are secondarily relevant, and a small value otherwise. Given a sample with expression annotations of $\hat{\mathbf{p}}_e = \{\hat{p}_{e0}, \hat{p}_{e1}, \dots, \hat{p}_{eE-1}\}$, it is intuitive to produce the pseudo AU labels by $\hat{\mathbf{p}}_e \widehat{\mathbf{W}}_{EA}$. However, the matrix merely defines three level correlations, which is very cursory. It is desirable to exploit finer-grained correlations so as to obtain more precise pseudo AU labels. In this work, we use a learnable correlation matrix \mathbf{W}_{EA} to generate the pseudo AU labels by

$$\hat{\mathbf{p}}_a = \hat{\mathbf{p}}_e \mathbf{W}_{EA}, \quad (1)$$

We apply the simple linear function to map the AU features $\{\mathbf{f}_1^a, \mathbf{f}_2^a, \dots, \mathbf{f}_A^a\}$ to the predicted AU labels

$$\mathbf{p}_a = \{\mathbf{w}_1 \mathbf{f}_1^a, \mathbf{w}_2 \mathbf{f}_2^a, \dots, \mathbf{w}_A \mathbf{f}_A^a\}, \quad (2)$$

where \mathbf{w}_a is a learnable weight vector. During training, we define a mean square error loss between the pseudo and the predicted labels, and a regularization loss between the learned and prior correlation matrices, formulated as

$$\mathcal{L}_{au} = \|\mathbf{p}_a - \hat{\mathbf{p}}_a\|_2^2 + \lambda \|\mathbf{W}_{EA} - \widehat{\mathbf{W}}_{EA}\|_2^2. \quad (3)$$

In this way, we can learn finer-grained correlation and simultaneously exploit prior correlations to promote generating more precise AU labels. λ is a balance parameter and it is set to 1.0 in the experiments.

C. Knowledge-Constrained AU Selection

In this subsection, we introduce the knowledge-constrained attention mechanism that learns to adaptively select useful AU features and fuses these features to enhance image representation. It computes a correlation coefficient for each AU, performs weighted average to obtain a merged AU feature, and concatenates it with expression feature to obtain the final image representation.

Specifically, we first fuse the expression features with each AU feature using the low-rank bilinear pooling [38] to compute a coefficient

$$\hat{w}_i = \mathbf{P}^T(\tanh(\mathbf{U}^T \mathbf{f}^e) \odot \tanh(\mathbf{V}^T \mathbf{f}_i^a) + \mathbf{b}) \quad (4)$$

that denotes the importance of AU i for expression recognition. Here, we use low-rank bilinear pooling [38] as it is effective for feature fusion. In the equation, $\tanh(\cdot)$ is the hyperbolic tangent function and \odot is the element-wise product operation. $U \in R^{d_e \times d}$, $V \in R^{d_a \times d}$, $P \in R^{d \times 1}$ are the learnable parameter matrixes. To make the coefficients easily comparable across different samples, we normalize the coefficients over all AUs using a softmax function

$$w_i = \frac{\hat{w}_i}{\sum_j \hat{w}_j}. \quad (5)$$

Then, we perform weighted average over all AUs to obtain the AU features

$$\mathbf{f}^a = \sum_i w_i \mathbf{f}_i^a. \quad (6)$$

Finally, we concatenate \mathbf{f}^a with \mathbf{f}^e for expression prediction

$$\mathbf{p}^e = f([\mathbf{f}^a, \mathbf{f}^e]). \quad (7)$$

D. Knowledge-Regularized Training Loss

As suggested in previous literatures [8], there exists strong co-occurrence dependencies among different AUs. In other words, some AUs co-occur frequently while some AUs are mutually exclusive. For example, the AU *inner brow raiser* always co-occurs with *outer brow raiser* as they are both controlled by the muscle group of *Occipito frontalis*. In contrast, the AU *lip corner puller* hardly co-exists with *lip corner depressor* as the corresponding controlled muscle group can not co-activate. It is expected and natural that the learned attentional coefficients should match such dependencies, and thus we introduce these prior dependencies as a regularization term during training.

Inspired by previous work [39], we consider the pair-wise dependencies that include positive and negative correlations to define the regularization term. Here, we consider the AU i exists if its attention coefficient is higher than 0.5. We denote it as i_1 if it exists and denote as i_0 otherwise. For the AU i and j with positive correlation, it is expected

$$\begin{aligned} p(i_1|j_1) &> p(i_0|j_1) \\ p(i_1|j_1) &> p(i_1|j_0) \end{aligned} \quad (8)$$

which can be transformed to the equivalent formulation

$$\begin{aligned} p(i_1, j_1) &> p(i_0, j_1) \\ p(i_1, j_1) &> p(i_1)p(j_1) \end{aligned} \quad (9)$$

Accordingly, we can define the regularization term that

constrains the positive correlation as

$$\begin{aligned} \ell_p &= \sum_{i,j \in S_p} \max(p(i_1)p(j_1) - p(i_1, j_1), 0) \\ &+ \sum_{i,j \in S_p} \max(p(i_1, j_0) - p(i_1, j_1), 0) \\ &+ \sum_{i,j \in S_p} \max(p(i_0, j_1) - p(i_1, j_1), 0). \end{aligned} \quad (10)$$

where S_p is the set of positive AU pairs. Similarly, the regularization term for negative correlation constraint can be defined as

$$\begin{aligned} \ell_n &= \sum_{i,j \in S_n} \max(p(i_1, j_1) - p(i_1)p(j_1), 0) \\ &+ \sum_{i,j \in S_n} \max(p(i_1, j_1) - p(i_1, j_0), 0) \\ &+ \sum_{i,j \in S_n} \max(p(i_1, j_1) - p(i_0, j_1), 0). \end{aligned} \quad (11)$$

where S_n is the set of negative AU pairs.

During training, we use the cross entropy loss, which is denoted by ℓ_c , to train the expression classifier, and thus the loss can be defined as

$$\mathcal{L}_e = \ell_c + \alpha(\ell_p + \ell_n), \quad (12)$$

where α is the balance parameter and it is set to 0.5 in the experiments.

E. Implementation Details

1) *Network architecture*: We select ResNet-101 [40] as the backbone network for feature extraction, which consists of four block layers. Given an input image of size 224×224 , we can obtain feature maps of size $56 \times 56 \times 256$ from first layer, feature maps of size $28 \times 28 \times 512$ from second layer, feature maps of size $14 \times 14 \times 1024$ from third layer and feature maps of size $7 \times 7 \times 2048$ from last layer. For holistic expression feature, we downsample the feature maps from the first, second, third layer of backbone network to the size of feature maps from last layer by max pooling, then concatenate these four feature maps, and perform global average pooling to obtain a 3840-dimensional vector. For AU feature, we inversely upsample the feature maps to the size of feature maps from previous layer by deconvolution, starting from feature maps from last layer to the end of second layer, and concatenate each upsampled feature maps with original feature maps, and upsample it again by deconvolution. At end, we obtain a feature map of size $56 \times 56 \times 256$, which is used for cropping feature from corresponding location to obtain the corresponding AU feature by using crop net.

In crop net, we first use MTCNN [41] to get facial landmarks of input image and crop the corresponding region for each AU on feature maps by using code of [13], and pass it through a convolutional layer and a fully connected layer, whose parameters are not shared for each AU, to obtain a 512-dimensional vector.

2) *Training details*: To obtain more stable experiment results, we adopt three-stages training process. In the first stage, we train the backbone and expression classifier with the cross-entropy loss using stochastic gradient descent(SGD) with an initial learning rate of 0.01, a momentum of 0.9, and a weight decay of 0. In the second stage, we fix the parameters of backbone, and train crop net and AU classifier with mean-square error loss and loss \mathcal{L}_{au} defined by the formula (3) using stochastic gradient descent(SGD) with an initial learning rate of 0.001, a momentum of 0.9, and a weight decay of 0. And in third stage, we fix the parameters of backbone, crop net and AU classifier, and train expression classifier with loss \mathcal{L}_e defined by the formula (12) using stochastic gradient descent(SGD) with an initial learning rate of 0.0001, a momentum of 0.9, and a weight decay of 0.

IV. EXPERIMENTS

A. Datasets

The existing datasets of facial expression can be divided into two main categories according to its collecting environment. Over quite a long period there only exist datasets collected in the lab-controlled conditions with limited size. Recently, some comparatively large-scale datasets that reflect real-world scenarios are released to promote the research. Meanwhile, the types of expressions are expanded with compound expressions that can be constructed by the combination of basic expression.

We chose two challenging in-the-wild datasets to evaluate the performance of our method. The Real-world Affective Face Database (RAF-DB) [4] and the Static Facial Expressions in the Wild (SFEW2.0) [5] datasets.

RAF-DB [4] contains 29,672 highly diverse facial images from thousands of individuals that were also collected from the Internet. With manually crowd-sourced annotation and reliable estimation, seven basic and eleven compound emotion labels are provided for the samples. Specifically, 15,339 images from the basic emotion set are divided into two groups (12,271 training samples and 3,068 testing samples) for evaluation.

SFEW2.0 [5] is an in-the-wild dataset collected from different films with spontaneous expressions, various head poses, age ranges, occlusions and illuminations. This dataset is divided into training, validation, and test sets, with 958, 436, and 372 samples, respectively.

B. Comparison with State-of-the-Arts

In this subsection, we present the performance comparisons with current state-of-the-art methods to evaluate the superiority of our proposed method.

1) *Performance on RAF-DB*: RAF-DB is a challenging in-the-wild dataset that is widely used for evaluating facial expression recognition. In this part, we compare our method with current state-of-the-art competitors, including Deep Neural Network Augmentation(DCNN-DA) [42], Weakly Supervised Local-Global Relation Network (WSLGRN) [43], Covariance Pooling (CP) [44], Compact Deep Learning Model(CompactDLM) [45], Feature Selection Network(FSN)

[46], Deep Locality-Preserving Learning (DLP-CNN) [4], and Multi-Region Ensemble CNN (MRE-CNN) [47].

We first present the accuracy of each basic expression and the average accuracy over all expressions in Table I. As shown, our method achieves the best performance compared with existing competitors, i.e., improving the average accuracy from 79.4% to 81.0%. In addition, our method obtains better accuracy for most basic expression, especially for those that are difficult to recognize. For example, current best accuracy for the expression ‘‘Fear’’ is 63.5%. Our method improves the accuracy to 68.9%, with a relative improvement of 8.50%. One possible reason is integrating information of facial AU can well capture local discriminative feature and help to distinguish uncertain and ambiguous expression.

Except for the basic expression, RAF-DB contains another sub-set in which each face image is annotated with a compound expression. A compound expression usually contains two or more basic expressions. For example, a person may be happy and simultaneously surprised. Obvious, this is an even more difficult task as it needs to recognize multiple expression patterns, which depends more on local discriminative feature mining. Here, we compare our method with BaseDCNN [4], Center Loss [4], and Deep Locality-Preserving Learning (DLP-CNN) [4]. As shown in Table II, our method outperforms current state-of-the-art competitors by a sizable margin, i.e., an improvement of 6.50% in average accuracy.

2) *Performance on SFEW2.0*: SFEW2.0 is an even more challenging dataset, and there are also some works that conduct experiments on this dataset. Here, we compare with our proposed method with the following works: Covariance Pooling (CP) [44], Deep Locality-Preserving Learning (DLP-CNN) [4], Identity-Aware Convolutional Neural Network (IA-CNN) [48], and Island Loss (IL) [49].

The accuracy of each basic expression and average accuracy overall all expressions are presented in Table III. Our method obtains an average accuracy of 52.8%, improving that of the previous best method by 1.7%. Similar to results on RAF, existing methods perform extremely poor for the expressions ‘‘Disgust’’ and ‘‘Fear’’, e.g., accuracies of 0.0% and 14.0% for these two expressions. Our methods improve the accuracies to 17.4% and 25.5%. These comparisons again demonstrate the superiority of our proposed method in ambiguous expression recognition.

C. Ablation Study

In this subsection, we conduct comprehensive ablation studies to discuss and analyze the contribution of each component and obtain a more thorough understanding of the framework.

1) *Analysis of AUE-CRL framework*: The core contribution of the proposed framework is mining useful local AU information to enhance image representation. To analyze the contribution of this framework, we compare AUE-CRL framework with the ResNet-101 baseline. The experiment is conducted on the RAF-DB dataset and the results are presented in Table IV. As shown, the average accuracy drops

Methods	Angry	Disgust	Fear	Happy	Neutral	Sad	Surprised	Ave. acc
DCNN-DA [42]	78.4	64.4	62.2	91.1	80.6	81.2	84.5	77.5
WSLGRN [43]	75.3	56.9	63.5	93.8	85.4	83.5	85.4	77.7
CP [44]	80.0	61.0	61.0	93.0	89.0	86.0	86.0	79.4
CompactDLM [45]	74.5	67.6	46.9	82.3	59.1	58.0	84.6	67.6
FSN [46]	72.8	46.9	56.8	90.5	76.9	81.6	81.8	72.5
DLP-CNN [4]	71.6	52.2	62.2	92.8	80.3	80.1	81.2	74.2
MRE-CNN [47]	84.0	57.5	60.8	88.8	80.2	79.9	86.0	76.7
Ours	80.5	67.6	68.9	94.1	85.8	83.6	86.4	81.0

TABLE I

PERFORMANCE OF OUR PROPOSED METHOD AND CURRENT EXISTING STATE-OF-THE-ART COMPETITORS FOR RECOGNIZING THE BASIC EXPRESSIONS ON THE RAF-DB DATASET.

Methods	BaseDCNN [4]	Center Loss [4]	DLP-CNN [4]	Ours
Ave. acc	40.2	40.0	44.6	51.1

TABLE II

PERFORMANCE OF OUR PROPOSED METHOD AND CURRENT EXISTING STATE-OF-THE-ART COMPETITORS FOR RECOGNIZING THE COMPOUND EXPRESSIONS ON THE RAF-DB DATASET.

from 81.0% to 79.9%. It is worth noting that, compared to AUE-CRL framework, the accuracy of baseline drops significantly on expression "Disgust", which proves AUE-CRL framework is effective to distinguish uncertain and ambiguous expression.

2) *Analysis of knowledge-guided AU representation learning*: As suggested above, we leverage the relationship of AU and expression to guide learning AU representation, which can help get rid of heavy AU annotations and guide learning domain-adaptive representation. To analyze the effect of the relationship of AU and expression, we remove the AU-expression regularization loss and use AU annotation of the BP4D dataset [37] to train AU classifiers. We conduct the comparison on the RAF-DB dataset and present the results on Table IV. Although this method adopts ground truth AU annotations, it performs inferior compared with our knowledge-guided AU representation learning, i.e., decreasing the accuracy from 81.0% to 80.2%.

Indeed, the images of BP4D cover merely 41 subjects and they are captured in the constrained lab environment. Thus, representation learned on such a dataset can hardly generalize to other environments. In contrast, our proposed knowledge-guided AU representation learning enables training on the target dataset and tends to learn domain-adaptive AU representation, leading to better performance.

3) *Analysis of knowledge-constrained AU selection*: In this work, we introduce a knowledge-constrained attention mechanism to adaptively select useful AU for expression representation enhancement. To analyze its contribution, we remove this component, simply perform average pooling over all action unit to obtain AU representation and concatenate it with expression feature for expression recognition. We find the average accuracy drops to 80.2% and accuracy drops significantly on expression "Disgust", which suggests that the

attention mechanism can help mine useful AU information to facilitate expression recognition and play a great role in distinguish uncertain and ambiguous expression.

To better select useful AUs, we introduce prior knowledge of the relationship among AUs as a constraint term during training. Here, we remove this constraint to analyze its contribution. As shown in Table IV, we find the average accuracy is 80.6% on the RAF-DB dataset, which is better than that without the attention mechanism but still worse than our method.

V. CONCLUSION

In this paper, we propose an AU-Expression Knowledge Constrained Representation Learning framework that exploits prior knowledge to help mining AU information to promote facial expression recognition. It first leverages relationships between AUs and expressions to guide learning domain-adaptive AU representation without any additional AU annotations. Then, it introduces an attentional mechanism to adaptively select useful AU representation under the constraint of the dependencies among AUs. We conduct an experiment on two in-the-wild datasets and show that our method outperforms current state-of-the-art competitors by a sizable margin.

REFERENCES

- [1] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*. IEEE, 2010, pp. 94–101.
- [2] M. Pantic, M. Valstar, R. Rademaker, and L. Maat, "Web-based database for facial expression analysis," in *2005 IEEE international conference on multimedia and Expo*. IEEE, 2005, pp. 5–pp.
- [3] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, "Coding facial expressions with gabor wavelets," in *Proceedings Third IEEE international conference on automatic face and gesture recognition*. IEEE, 1998, pp. 200–205.
- [4] S. Li, W. Deng, and J. Du, "Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [5] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon, "Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark," in *IEEE International Conference on Computer Vision Workshops*, 2011.
- [6] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

Methods	Angry	Disgust	Fear	Happy	Neutral	Sad	Surprised	Ave. acc
CP [44]	66.0	0.0	14.0	90.0	86.0	66.0	29.0	50.1
DLP-CNN [4]	-	-	-	-	-	-	-	51.1
IA-CNN [48]	70.7	0.0	8.9	70.4	60.3	58.8	28.9	42.6
IL [49]	61.0	0.0	6.4	89.0	66.2	48.0	33.3	43.4
Ours	75.3	17.4	25.5	86.3	72.1	50.7	42.1	52.8

TABLE III

PERFORMANCE OF OUR PROPOSED METHOD AND CURRENT EXISTING STATE-OF-THE-ART COMPETITORS RECOGNIZING THE BASIC EXPRESSIONS ON THE SFEW2.0 DATASET. - DENOTES CORRESPONDING RESULT IS NOT PROVIDED.

Methods	Surprised	Fear	Disgust	Happy	Sad	Angry	Neutral	Ave. acc
Baseline	85.5	68.9	60.1	94.1	85.3	81.8	83.3	79.9
Ours w/o KGAURL	88.3	63.5	64.2	92.7	85.7	79.9	87.3	80.2
Ours w/o Attention	88.0	67.6	60.1	94.1	84.0	82.5	85.0	80.2
Ours w/o AU-Independent	86.1	67.6	65.5	94.2	84.6	79.9	85.9	80.6
Ours	86.4	68.9	67.6	94.1	83.6	80.5	85.8	81.0

TABLE IV

PERFORMANCE OF OUR METHOD (OURS), OUR METHOD WITHOUT KNOWLEDGE-GUIDED AU REPRESENTATION LEARNING (OURS W/O KGAURL), OUR METHOD WITHOUT ATTENTION MECHANISM FOR AU REPRESENTAION SELECTION (OURS W/O ATTENTION), OUR METHOD WITHOUT AU-INDENPENT CONSTRAIN (OURS W/O AU-INDENPENDENT), AND THE BASELINE RESNET 101 (BASELINE) ON THE RAF-DB DATASET.

- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [8] E. Friesen and P. Ekman, "Facial action coding system: a technique for the measurement of facial movement," *Palo Alto*, vol. 3, 1978.
- [9] Y.-I. Tian, T. Kanade, and J. F. Cohn, "Recognizing action units for facial expression analysis," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 23, no. 2, pp. 97–115, 2001.
- [10] M. F. Valstar, M. Mehu, B. Jiang, M. Pantic, and K. Scherer, "Meta-analysis of the first facial expression recognition challenge," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 42, no. 4, pp. 966–979, Aug 2012.
- [11] M. Dahmane and J. Meunier, "Emotion recognition using dynamic grid-based hog features," in *Face and Gesture 2011*, March 2011, pp. 884–888.
- [12] Y. Hu, Z. Zeng, L. Yin, X. Wei, X. Zhou, and T. S. Huang, "Multi-view facial expression recognition," in *2008 8th IEEE International Conference on Automatic Face Gesture Recognition*, Sep. 2008, pp. 1–6.
- [13] G. Li, X. Zhu, Y. Zeng, Q. Wang, and L. Lin, "Semantic relationships guided representation learning for facial action unit recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8594–8601.
- [14] A. Mollahosseini, B. Hasani, and M. H. Mahoor, "Affectnet: A database for facial expression, valence, and arousal computing in the wild," *IEEE Transactions on Affective Computing*, 2017.
- [15] P. Liu, S. Han, Z. Meng, and Y. Tong, "Facial expression recognition via a boosted deep belief network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1805–1812.
- [16] A. Mollahosseini, D. Chan, and M. H. Mahoor, "Going deeper in facial expression recognition using deep neural networks," in *2016 IEEE Winter conference on applications of computer vision (WACV)*. IEEE, 2016, pp. 1–10.
- [17] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [18] M. Liu, S. Li, S. Shan, R. Wang, and X. Chen, "Deeply learning deformable facial action parts model for dynamic expression analysis," D. Cremers, I. Reid, H. Saito, and M.-H. Yang, Eds., 2015, pp. 143–157.
- [19] Z. Yu and C. Zhang, "Image based static facial expression recognition with multiple deep network learning," in *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, ser. ICMI '15, 2015, pp. 435–442.
- [20] Y. Xie, T. Chen, T. Pu, H. Wu, and L. Lin, "Adversarial graph representation adaptation for cross-domain facial expression recognition," in *Proceedings of the 28th ACM international conference on Multimedia*, 2020, pp. 1255–1264.
- [21] T. Chen, T. Pu, Y. Xie, H. Wu, L. Liu, and L. Lin, "Cross-domain facial expression recognition: A unified evaluation benchmark and adversarial graph learning," 2020.
- [22] T. Chen, L. Lin, X. Hui, R. Chen, and H. Wu, "Knowledge-guided multi-label few-shot learning for general image recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [23] T. Chen, M. Xu, X. Hui, H. Wu, and L. Lin, "Learning semantic-specific graph representation for multi-label image recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 522–531.
- [24] T. Chen, W. Yu, R. Chen, and L. Lin, "Knowledge-embedded routing network for scene graph generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6163–6171.
- [25] Z. Wang, T. Chen, J. Ren, W. Yu, H. Cheng, and L. Lin, "Deep reasoning with knowledge graph for social relationship understanding," in *Proc. of International Joint Conference on Artificial Intelligence*, 2018, pp. 2021–2028.
- [26] Y. Li, D. Tarlow, M. Brockschmidt, and R. S. Zemel, "Gated graph sequence neural networks," in *International Conference on Learning Representations*, 2016.
- [27] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.
- [28] P. Ekman and W. Friesen, "Facial action coding system: A technique for the measurement of facial movement," *Consulting Psychologists Press*, 1978.
- [29] A. Savran, B. Sankur, and M. T. Bilge, "Regression-based intensity estimation of facial action units," *Image and Vision Computing*, vol. 30, no. 10, pp. 774–784, 2012.
- [30] M. H. Mahoor, S. Cadavid, D. S. Messinger, and J. F. Cohn, "A framework for automated measurement of the intensity of non-posed facial action units," in *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2009, pp. 74–80.
- [31] J. C. Batista, V. Albiero, O. R. Bellon, and L. Silva, "Aumpnet: simultaneous action units detection and intensity estimation on multipose facial images using a single convolutional neural network," in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 866–871.
- [32] Y. Zhou, J. Pi, and B. E. Shi, "Pose-independent facial action unit intensity regression based on multi-task deep transfer learning," in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 872–877.

- [33] R. Walecki, V. Pavlovic, B. Schuller, M. Pantic, *et al.*, “Deep structured learning for facial action unit intensity estimation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3405–3414.
- [34] A. Gudi, H. E. Tasli, T. M. Den Uyl, and A. Maroulis, “Deep learning based faces action unit occurrence and intensity estimation,” in *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, vol. 6. IEEE, 2015, pp. 1–5.
- [35] T. Chen, L. Lin, W. Zuo, X. Luo, and L. Zhang, “Learning a wavelet-like auto-encoder to accelerate deep neural networks,” in *Proc. of AAAI Conference on Artificial Intelligence*, 2018, pp. 6722–6729.
- [36] S. M. Mavadati, M. H. Mahoor, K. Bartlett, P. Trinh, and J. F. Cohn, “Disfa: A spontaneous facial action intensity database,” *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 151–160, 2013.
- [37] X. Zhang, L. Yin, J. F. Cohn, S. Canavan, M. Reale, A. Horowitz, P. Liu, and J. M. Girard, “Bp4d-spontaneous: a high-resolution spontaneous 3d dynamic facial expression database,” *Image and Vision Computing*, vol. 32, no. 10, pp. 692–706, 2014.
- [38] J. H. Kim, K. W. On, J. Kim, J. W. Ha, and B. T. Zhang, “Hadamard product for low-rank bilinear pooling,” 2016.
- [39] Y. Zhang, W. Dong, B.-G. Hu, and Q. Ji, “Classifier learning with prior probabilities for facial action unit recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5108–5116.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [41] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, “Joint face detection and alignment using multi-task cascaded convolutional networks,” *CoRR*, vol. abs/1604.02878, 2016.
- [42] D. Kollias, S. Cheng, E. Ververas, I. Kotsia, and S. Zafeiriou, “Deep neural network augmentation: Generating faces for affect analysis,” *International Journal of Computer Vision*, 02 2020.
- [43] H. Zhang, W. Su, and J. Y. and Zengfu Wang, “Weakly supervised local-global relation network for facial expression recognition,” in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, 2020, pp. 1040–1046.
- [44] D. Acharya, Z. Huang, D. Pani Paudel, and L. Van Gool, “Covariance pooling for facial expression recognition,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.
- [45] C. Kuo, S. Lai, and M. Sarkis, “A compact deep learning model for robust facial expression recognition,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 2202–22028.
- [46] S. Zhao, H. Cai, H. Liu, J. Zhang, and S. Chen, “Feature selection mechanism in cnns for facial expression recognition,” in *BMVC*, 2018, p. 317.
- [47] Y. Fan, J. C. Lam, and V. O. Li, “Multi-region ensemble convolutional neural network for facial expression recognition,” in *International Conference on Artificial Neural Networks*. Springer, 2018, pp. 84–94.
- [48] Z. Meng, P. Liu, J. Cai, S. Han, and Y. Tong, “Identity-aware convolutional neural network for facial expression recognition,” in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 558–565.
- [49] J. Cai, Z. Meng, A. S. Khan, Z. Li, J. O’Reilly, and Y. Tong, “Island loss for learning discriminative features in facial expression recognition,” in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 302–309.