

# Region Attention Networks for Pose and Occlusion Robust Facial Expression Recognition

Kai Wang, Xiaojiang Peng, Jianfei Yang, Debin Meng, and Yu Qiao, *Senior Member, IEEE*

**Abstract**—Occlusion and pose variations, which can change facial appearance significantly, are two major obstacles for automatic Facial Expression Recognition (FER). Though automatic FER has made substantial progresses in the past few decades, occlusion-robust and pose-invariant issues of FER have received relatively less attention, especially in real-world scenarios. This paper addresses the real-world pose and occlusion robust FER problem in the following aspects. First, to stimulate the research of FER under real-world occlusions and variant poses, we annotate several in-the-wild FER datasets with pose and occlusion attributes for the community. Second, we propose a novel Region Attention Network (RAN), to adaptively capture the importance of facial regions for occlusion and pose variant FER. The RAN aggregates and embeds varied number of region features produced by a backbone convolutional neural network into a compact fixed-length representation. Last, inspired by the fact that facial expressions are mainly defined by facial action units, we propose a region biased loss to encourage high attention weights for the most important regions. We validate our RAN and region biased loss on both our built test datasets and four popular datasets: FERPlus, AffectNet, RAF-DB, and SFEW. Extensive experiments show that our RAN and region biased loss largely improve the performance of FER with occlusion and variant pose. Our method also achieves state-of-the-art results on FERPlus, AffectNet, RAF-DB, and SFEW. Code and the collected test data will be publicly available.

**Index Terms**—Facial expression recognition, occlusion-robust and pose-invariant, region attention network, deep convolutional neural networks.

## I. INTRODUCTION

Facial expressions play important roles in daily human-human communication. Automatic facial expression analysis is an important area of artificial intelligence. Due to its potential applications in various fields, such as intelligent tutoring systems, service robots, driver fatigue monitoring, Facial Expression Recognition (FER) has attracted increasing attention in the computer vision community recently [28], [5], [21], [48], [36], [17]. The main challenges of FER come from illumination variation, occlusions, variant poses, identity bias, insufficient qualitative data, etc.

Occlusions and variant poses are two major problems in the field of face analysis since they lead to significant change of facial appearance. These issues have received wide interest in face identity recognition [19], [38], however, less attention has

been paid in real-world FER partly due to the lack of a facial expression dataset with occlusion and pose annotations.

Earlier works mainly investigate the effects of occlusion for FER systems with partial artificially-masked faces collected in a controlled laboratory environment. Boucher and Ekman [7] investigate facial parts to understand which are most important regions for human perception by occluding key parts. Bourel *et al.* [8] present the first FER system under the presence of occlusion by recovering geometric facial points. Kotsia *et al.* [32] present a comprehensive analysis on occluded FER based on Gabor features and human observers, and find that an occluded mouth degrades FER more than occluded eyes on JAFFE [44] and CK [30]. Sparse representation classifier (SRC) is widely used for artificially-occluded FER in 2010s [12], [71], [13]. Subsequently, a number of works handle FER with sub-region based features and fusion schemes [24], [69], [70], which detect the occlusion regions first and then remove their local features. With the popularity of data-driven deep learning techniques, several recent efforts on FER have been made on the collection of large-scale datasets [21], [48], [36], and many works [57], [29], [40], [6], [1], [35] exploit deep convolutional neural networks (CNN) to improve the performance of FER.

We argue that explicitly removing occlusion regions is not practical since real-world occlusion is difficult to detect in itself. Directly using CNN on whole face images ignores the characteristics of occlusion and variant pose. In practices, occlusion and pose variations can lead to unseen regions of input faces, which bring difficulties for face alignment and harm the feature extraction process. Contrasted with these difficulties, human have the remarkable ability to understand facial expressions under challenging conditions. Psychological studies indicated that human can effectively exploit both local regions and holistic faces to perceive the semantics delivered through incomplete faces [64]. Inspired by these facts, this paper proposes a region based deep attention architecture for pose and occlusion robust FER, which adaptively integrates visual clues from regions and whole faces. Specifically, we addresses the real-world pose and occlusion robust FER problem in the following aspects.

First, to investigate the occlusion and pose variant FER problem, we build six real-world test datasets from FERPlus and AffectNet, namely Occlusion-FERPlus, Pose-FERPlus, Occlusion-AffectNet, and Pose-AffectNet, Occlusion-RAF-DB, and Pose-RAF-DB. The occlusion test datasets are manually annotated with occlusion types of wearing mask/glasses, objects in left/right, objects in upper face, objects in bottom face. The pose-variant test datasets are automatically labeled by a recent head pose estimation toolbox [3]. We observe that

Kai Wang and Xiaojiang Peng are equally-contributed authors.

Kai Wang, Xiaojiang Peng, Debin Meng and Yu Qiao are with Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China.

Jianfei Yang is with School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore

Corresponding Author: Yu Qiao(yu.qiao@siat.ac.cn)

the performance of existing CNN methods degrade significantly in occlusion and pose-variant environments.

Second, we propose the Region Attention Network (RAN), to capture the importance of facial regions for occlusion and pose robust FER. The RAN is comprised of a feature extraction module, a self-attention module, and a relation attention module. The later two modules aim to learn coarse attention weights and refine them with global context, respectively. Given a number of facial regions, our RAN learns attention weights for each regions in an end-to-end manner, and aggregates their CNN-based features into a compact fixed-length representation. Besides, the RAN model has two auxiliary effects on the face images. On one hand, *cropping regions can enlarge* the training data which is important for those insufficient challenging samples. On the other hand, *rescaling the regions to the size of original images highlights fine-grain facial features*. Extensive experiments indicate that our RAN significantly improves the performance of FER in occlusion and pose variant conditions.

Third, since facial expressions are mainly defined by multiple facial action units [7], we propose a Region Biased Loss (RB-Loss) to encourage a high attention weight for the most important region. Our RB-Loss resorts a simple constraint on the RAN that the maximum attention weight of facial regions should be larger than the one of the original face image. Experiments show that the RB-Loss further improves FER slightly without additional computation cost. Our FER solution achieves state-of-the-art results on FERPlus, AffectNet, RAF-DB, and SFEW with accuracies of **89.16%**, **59.5%**, **86.9%**, and **56.4%**, respectively.

## II. RELATED WORK

In this section, we mainly present related works on normal FER problem, the occlusion and pose variant FER problem, and attention mechanism.

**Facial Expression Recognition.** Generally, a FER system mainly consists of three stages, namely face detection, feature extraction, and expression recognition. In face detection, several face detectors like MTCNN [67] and Dlib [3]) are used to locate faces in complex scenes. The detected faces can be further aligned alternatively. For feature extraction, various methods are designed to capture facial geometry and appearance features caused by facial expressions. According to the feature type, they can be grouped into engineered features and learning-based features. For the engineered features, they can be further divided into texture-based local features, geometry-based global features, and hybrid features. The texture-based features mainly include SIFT [49], HOG [14], Histograms of LBP [52], Gabor wavelet coefficients [39], etc. The geometry-based features are mainly based on the landmark points around noses, eyes, and mouths. Combining two or more of the engineered features refers to the hybrid feature extraction, which can further enrich the representation. For the learned features, Fasel [22] finds that a shallow CNN is robust to face poses and scales. Tang [57] and Kahou *et al.* [29] utilize deep CNNs for feature extraction, and win the FER2013 and EmotiW2013 challenge, respectively. Liu *et al.* [40] propose

a Facial Action Units based CNN architecture for expression recognition. After feature extraction, the next stage is to feed the features into a supervised classifier such as Support Vector Machines (SVMs), softmax layer, and logistic regression to assign expression categories.

To avoid overfitting on small facial expression datasets, many recent studies [45], [2], [56], [33], [72], [20] utilize face recognition datasets to pre-train a network, and then fine-tune it on target expression datasets. Levi and Hassner [33] leverage the CASIA-WebFace[63] face recognition dataset to pretrain four different VGGNet[53] and GoogleNet[55]. Zhao *et al.* [72] propose a Peak Gradient Suppression (PGS) scheme for training and also pretrain their models on CASIA-WebFace. Ding *et al.* [20] propose a FaceNet2ExpNet framework which jointly trains FER task and face recognition task. Albanie *et al.* use the VGGFace model (face recognition model) and fine-tune it on FERPlus with soft probabilities. Meng *et al.* evaluate different face recognition model architectures and used face recognition datasets for facial expression.

**FER in Occlusion and Pose Variant Condition.** Occlusion and variant pose usually occur in real-world scenarios as facial regions can be easily occluded by sunglasses, a hat, a scarf, etc. Partial occlusion can be divided into two types according to whether the real object causes occlusion: one is artificial occlusion, and the other is real-life occlusion. Few attempts have been made on the real-world occlusion FER problem. Kotsia *et al.* [32] demonstrate how artificial partial occlusion affects the FER, and discuss how to deal with it. Liu *et al.* [41] propose a novel FER method to address partial occlusion problem based on Gabor multi-orientation features fusion and local Gabor binary pattern histogram sequence (LGBPHS). Cotter *et al.* [12], [13] propose to use sparse representation classifier for partial occlusion FER. The latest related work [37] designs a patch-based attention network for occlusion aware FER. The patches are cropped from the area of eyes, nose, mouth and so on. The selected 24 patches are fed into an attention network which is near to the self-attention module in our work. Our work differs from [37] in that i) we crop relative large regions instead of small fixed parts by considering that the facial expression is connected to multiple AUs, and ii) we refine the attention weights with a relation-attention module and region bias loss function. As for the pose variant FER problem, Rudovic *et al.* [51] propose the Coupled Scaled Gaussian Process Regression (CSGPR) model for head-pose normalization. Different from existing methods, we address both occlusion and pose variant FER problems in an end-to-end manner with an elaborately-designed region attention network architecture and collected test datasets.

**Attention Networks.** Attention mechanisms are firstly developed on the basis of reinforcement algorithm. Mnih *et al.* [47] use the attention on the RNN model for image classification, and then it is successfully utilized for machine translation tasks. Bahdanau *et al.* [4] use an attention-like mechanism to simultaneously translate and align the source languages, and their work is the first attempt to apply attention mechanism to machine translation. Afterwards, many self-attention models are proposed for different tasks, such as LSTM for machine

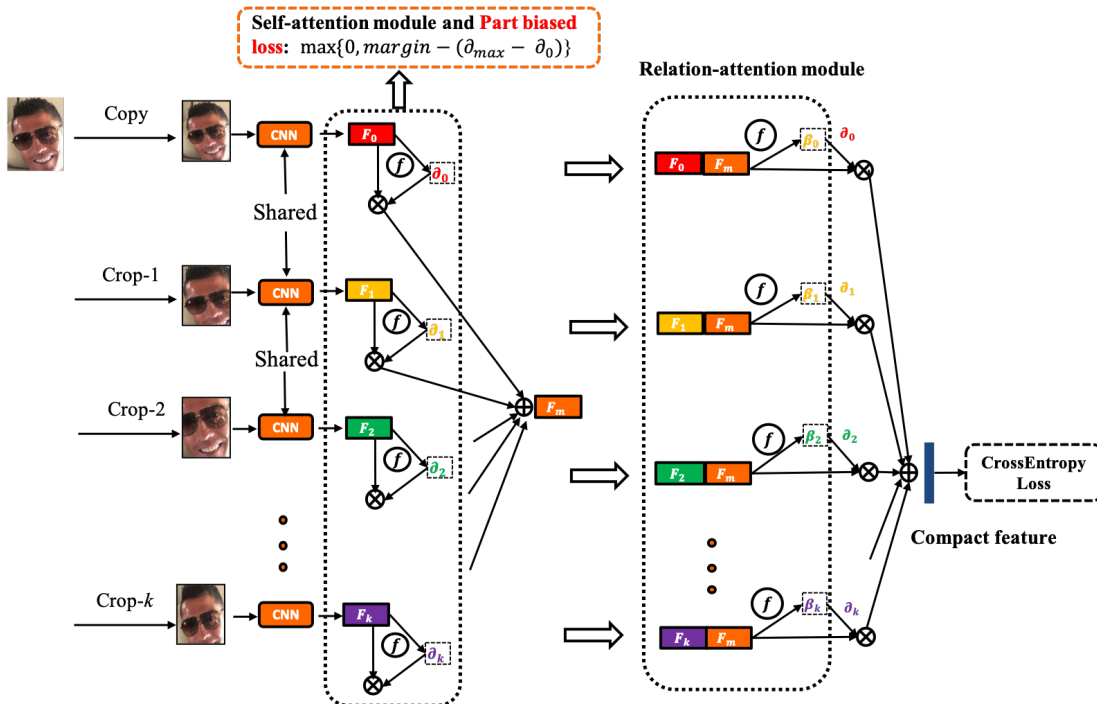


Fig. 1: The framework of our RAN. A face image is cropped into several regions, and these regions are fed into a backbone CNN for feature extraction. The self-attention and relation-attention module are then used to obtain compact face representation.  $\sigma$  denotes the sigmoid function.

reading [11], multi-head attention for machine translation [58] and attention clusters for video classification [42]. Yu et al. [59] propose an attention network for face detection, which highlights the face regions in anchor generation step. Perhaps the most similar work to ours is the Neural Aggregation Network (NAN) proposed by Yang et al. [62]. NAN uses a cascade attention mechanism to aggregate face features of a video or set into a compact video representation. Our work differs from NAN by that self-attention and relation-attention module is used in RAN to aggregate facial region features for FER in static images, and a region biased loss is introduced to enhance region weights.

### III. METHODOLOGY

In this section, we first give an overview of our proposed region attention networks (RAN), and then detail each module and the region biased loss in RAN. We then present the region generation strategies and finally describe the collected occlusion and pose variant FER dataset.

#### A. Overview

As mentioned in Sec. I, several early works try to detect the occlusion regions and then remove the region features to address the facial expression recognition with regional occlusion. Along with this idea, we aim to automatically reduce or eliminate the effect of occlusion and irrelevant regions with an end-to-end deep architecture.

Considering both large pose and occlusion issues in facial expression recognition, we propose a Region Attention Network (RAN) to alleviate the degradation of naive face based

CNN models. The proposed RAN can adaptively capture the importance of facial region information, and make a reasonable trade-off between region and global features. The pipeline of our RAN is illustrated in Figure 1. It mainly consists of three modules, namely region cropping and feature extraction module, self-attention module, and relation-attention module. Given a face image (after face detection), we first crop it into a number of regions with fixed position cropping or random cropping. We will compare these strategies in experiments. These regions along with the original face region are then fed into a backbone CNN model for region feature extraction. Subsequently, the self-attention module assigns an attention weight for each region using a fully-connected (FC) layer and the sigmoid function. An alternative region biased loss (RB-Loss) is further introduced to regularize the attention weights and enhance the most valuable region in self-attention module. We aggregate these region features to a global representation ( $F_m$  in Figure 1). Then the relation-attention module uses a similar attention mechanism on the concatenation of individual region feature and global representation to further capture content-aware attention weights. Finally, we leverage the weighted region feature and the global representation to predict the expressions.

#### B. Region Attention Networks

As shown in Figure 1, the proposed RAN mainly consists of two stages. The first stage is to coarsely calculate the importance of each region by a FC layer conducted on its own feature, which is called self-attention module. The second stage seeks to find more accurate attention weights

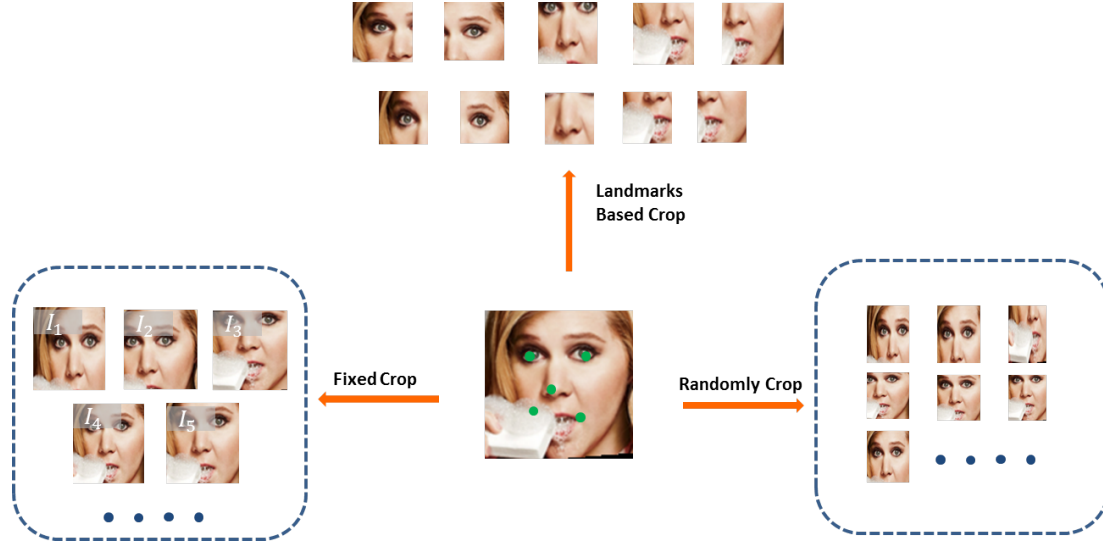


Fig. 2: An example of our region generation methods. Left: fixed position cropping. Right: random cropping. Upper: landmark-based cropping.

by modeling the relation between the region features and the aggregated content representation from the first stage, which is called relation-attention module.

Formally, we denote a face image as  $I$ , its duplicate as  $I_0$ , and its crops as  $I_1, I_2, \dots, I_k$ , and the backbone CNN as  $r(\cdot; \theta)$ . The feature set  $X$  of  $I$  is defined by:

$$X = [F_0, F_1, \dots, F_k] = [r(I_0; \theta), r(I_1; \theta), \dots, r(I_k; \theta)], \quad (1)$$

where  $\theta$  is the parameter of backbone CNN.

**Self-attention module.** With these region features, the self-attention module applies a FC layer and a sigmoid function to estimate coarse attention weights. Mathematically, the attention weight of the  $i$ -th region is defined by:

$$\mu_i = f(F_i^\top \mathbf{q}^0), \quad (2)$$

where  $\mathbf{q}^0$  is the parameter of FC,  $f$  denotes the sigmoid function. In this stage, we summarize all the region features with their attention weights into a global representation  $F_m$  as follows,

$$F_m = \frac{1}{\sum_{i=0}^n \mu_i} \sum_{i=0}^n \mu_i F_i. \quad (3)$$

$F_m$  is a compact representation and can be used as the final input of classifier. We compare the self-attention aggregation to the straightforward average pooling and concatenation (with fixed number of crops) in Sec. IV.

**Relation-attention module.** The self-attention module learns weights with individual features and non-linear mapping, which is rather coarse. Since the aggregated representation  $F_m$  inherently represents the contents of all the facial regions, the attention weights can be further refined by modeling the relation between region features and this global representation  $F_m$ .

Inspired by the global attention in neural machine translation [43] and the relation-Net in low-shot learning [61], we use the sample concatenation and another FC layer to estimate new

attention weights for region features. The new attention weight of the  $i$ -th region in relation-attention module is formulated as,

$$\nu_i = f([F_i : F_m]^\top \mathbf{q}^1), \quad (4)$$

where  $\mathbf{q}^1$  is the parameter of FC, and  $f$  denotes the sigmoid function. In this stage, we aggregate all the region information along with the coarse global representation from self-attention into a new compact feature as,

$$P_{RAN} = \frac{1}{\sum_{i=0}^n \mu_i \nu_i} \sum_{i=0}^n \mu_i \nu_i [F_i : F_m]. \quad (5)$$

$P_{RAN}$  is used as the final representation of the proposed RAN method.

**Region Biased Loss.** Inspired by the observation that different facial expressions are mainly defined by different facial regions [7], we make a straightforward constraint on the attention weights of self-attention, *i.e.* region biased loss (RB-Loss). This constraint enforces that one of the attention weights from facial crops should be larger than the original face image with a margin. For example, the Crop-2 in Figure 1 can be more discriminative than the original one. Formally, the RB-Loss is defined as,

$$\mathcal{L}_{RB} = \max\{0, \alpha - (\mu_{max} - \mu_0)\}, \quad (6)$$

where  $\alpha$  is a hyper-parameter served as a margin,  $\mu_0$  is the attention weight of the copy face image,  $\mu_{max}$  denotes the maximum weight of all facial crops.

In training, the classification loss is jointly optimized with the region biased loss. The proposed RB-Loss enhances the effect of region attention and encourages RAN to obtain superior weights of region and global representations. In fact, the RB-Loss can be also added to the relation-attention module. However, since the features of the relation-attention module already include holistic information, we experimentally find there is no gain by adding RB-Loss on the relation attention module.



Fig. 3: Some examples of our collected occlusion and pose variant test datasets. The left color images are from the test set of AffectNet, and the right gray images are from the test set of FERPlus.

### C. Region Generation

Cropping multiple regions is a fundamental step of our RAN. Too large regions lead to the reduced diversity of features and degrade to the case of many duplicates of the original face. Too small regions lead to insufficient discrimination ability of region features. In this paper, we evaluate three kinds of region generation schemes for our region attention networks, namely fixed position cropping, random cropping, and landmark-based cropping which are depicted in Figure 2.

**Fixed position cropping.** Since the face image can be well aligned by the recent advanced face alignment methods, a simple region generation scheme is to crop regions in fixed positions with fixed scales. Specifically, we crop five regions. Three of them are the top-left, top-right and center-down face regions, which have fixed size of 0.75 scale ratio of the original face. The other two regions are similar to those used in the smile-classification task of [68]. Here we crop the center regions with sizes of 0.9, and 0.85 scale ratio of the original face. All the crops are resized to have the same input size of the backbone CNN.

**Random cropping.** In deep face recognition, the DeepID method uses 200 random crops for each face image to enhance its performance [54]. For random cropping in our approach, we randomly crop  $N$  regions with random sizes ranged from 0.7 to 0.95 scale ratio of the original face.

**Landmark-based cropping.** Given facial landmarks, a straightforward method is to crop regions surrounding them, which is also used in [37]. Here we use MTCNN to detect five facial landmarks (i.e. left eye, right eye, nose, left mouth corner, and right mouth corner), and use them to crop five regions. Specifically, according to each facial landmark, we use a radius  $r$  to crop regions and remove the regions which are out of the original image.

### D. Occlusion and Pose Variant Dataset

Though our proposed RAN can be used for FER in any conditions, we focus on the real-world occlusion and pose variation problems. To the best of our knowledge, there is only a small real-world occlusion test dataset released in [37] very recently, and there is no publicly available facial expression dataset that addresses both occlusion and pose annotations. To examine our method under real-world scenario, we build six test datasets from the existing large-scale FER datasets. From the test set of FERPlus [5], the validation set of AffectNet [48], and the test set of RAF-DB [34], we collect the Occlusion-FERPlus, Pose-FERPlus, Occlusion-AffectNet, Pose-AffectNet, Occlusion-RAF-DB, and Pose-RAF-DB for testing. The test set will be available at <https://github.com/kaiwang960112/Challenge-condition-FER-dataset>. These real-world test sets are annotated with different occlusion types and different pose degrees. Some examples are illustrated in Figure 3.

For the pose variant test sets, we use the popular OpenFace toolbox [3] to estimate the Euler Angle in pitch, yaw, roll directions. Since the roll angle is in-plane which can be eliminated by face alignment, we only consider the pose in pitch and yaw directions. Those faces with pitch or yaw angle larger than  $30^\circ$  are collected to Pose-FERPlus, Pose-AffectNet, and Pose-RAF-DB.

For the occlusion test sets, we first define several occlusion types, namely wearing mask, wearing glasses, objects in left/right, objects in upper face and objects in bottom face, non-occlusion. Then we manually assign these categories to the test sets of FERPlus, AffectNet, and RAF-DB. Images with at least one type of occlusion are selected as the occlusion test sets.

We present the statistics of our collected test sets in Table I. Among all the occlusion types on FERPlus, AffectNet, and RAF-DB, the upper occlusion has the smallest samples.

TABLE I: Statistics of collected test datasets.

	Occlusion			Pose(pitch/yaw)		
	upper	bottom	left/right	glasses/mask	>30	>45
FERPlus	70	138	213	184	1171	634
AffectNet	84	183	128	288	1949	985
RAF-DB	126	151	160	298	1248	558

The total numbers of occlusion samples in FERPlus (test) , AffectNet (validation), and RAF-DB (test) are respectively 605, 682 , and 735, which are 16.86%, 17.05%, and 23.9% of their original sets. For the variant pose issue, about one-third of FERPlus (test), about two-fifths of RAF-DB, and about half of AffectNet (validation) have poses larger than 30 degrees (in pitch or yaw).

#### IV. EXPERIMENTS

In this section, we first describe the used datasets and our implementation details. We then present our collected occlusion and pose variant test datasets and evaluate our proposed RAN on them. We further explore each components of RAN on FERPlus [5], AffectNet [48] , and SFEW [15]. Finally, we compare our method to the state-of-the-art approaches.

##### A. Datasets

To evaluate our method, we use four popular in-the-wild facial expression datasets, namely FERPlus [5], AffectNet [48], RAF-DB [34], and SFEW [15]. These datasets cover different scales of face images and the challenging conditions. Besides, we also build occlusion and pose variant test datasets from FERPlus, AffectNet, and RAF-DB.

**FERPlus** [5]. The FERPlus is extent from FER2013 [23] introduced during the ICML 2013 Challenges in Representation Learning. It is a large-scale and real-world datasets collected by the Google search engine, and consists of 28,709 training images, 3,589 validation images and 3,589 test images. All face images in the dataset are aligned and resized to  $48 \times 48$ . The main difference between FER2013 and FERPlus is the annotation. FER2013 is annotated with seven expression labels (neutral, happiness, surprise, sadness, anger, disgust, fear) by one tagger, while FERPlus adds a *contempt* label and is annotated by 10 labels. In [5], the authors evaluate several training schemes, such as one-hot label (majority voting) and label distribution with cross-entropy loss. We mainly report the overall accuracy on the test set with supervision of majority voting and label distribution.

**AffectNet** [48]. The AffectNet is by far the largest dataset that provides both categorical and Valence-Arousal annotations. The dataset contains more than one million images from Internet by querying expression-related keywords in three search engines, of which 450,000 images are manually annotated with eight basic expression labels as FERPlus. AffectNet has an imbalanced test set, a balanced validation set, and an imbalanced training set. We mainly report accuracy on the validation set where each category contains 500 samples.

**SFEW** [15]. The Static Facial Expressions in the Wild (SFEW) dataset is built by selecting frames from AFEW [16], which covers unconstrained facial expressions, varied head

TABLE II: Performance comparison between the proposed RAN and baseline method with occlusion and variant pose conditions.

	FERPlus	Occlusion	Pose(30)	Pose(45)
Baseline		73.33	78.11	75.50
RAN (w RB-Loss)		<b>83.63</b>	<b>82.23</b>	<b>80.40</b>
	AffectNet	Occlusion	Pose(30)	Pose(45)
Baseline		49.48	50.10	48.50
RAN (w RB-Loss)		<b>58.50</b>	<b>53.90</b>	<b>53.19</b>
	RAF-DB	Occlusion	Pose(30)	Pose(45)
Baseline		80.19	84.04	83.15
RAN (w RB-Loss)		<b>82.72</b>	<b>86.74</b>	<b>85.20</b>

poses, large age range, occlusions, varied focus, different resolution of the face and real-world illumination. We use the newest version of SFEW in [18] where it has been divided into three sets: train (958 images), validation (436 images), and test (372 images). Each image is labeled with one of the seven expressions including angry, disgust, fear, happy, sad, surprise, and neutral by two independent labelers. We mainly report our performance on the validation set.

**RAF-DB.** RAF-DB [34] contains 30,000 facial images annotated with basic or compound expressions by 40 trained human coders. In our experiment, only images with basic emotions were used, including 12,271 images as training data and 3,068 images as test data.

##### B. Implementation details

In all the following experiments, we use the CNN detector and the ERT [31] based face alignment method in Dlib toolbox<sup>1</sup> to crop and align faces, and then resize them to the size of  $224 \times 224$ . We implement our methods with Pytorch toolbox<sup>2</sup>. For the backbone CNN, we mainly use the ResNet-18 [25] and VGG16 [50]. The ResNet-18 is pre-trained on MS-Celeb-1M face recognition dataset and VGG16 is downloaded from website<sup>3</sup>. The last pooling layer of ResNet-18, and the first FC feature of VGG16 is used for facial representation. In training phase of fixed cropping, we use all the five regions along with original face for each face image (i.e.  $k = 5$  in Figure 1). For training with random cropping, we replace the fixed five regions with randomly cropped ones. When jointly training with RB-Loss and Cross-Entropy loss, the default loss weight ratio is 1:1. On all datasets, the learning rate is initialized as 0.01, and divided by 10 after 15 epochs and 30 epochs. We stop training in 40 epochs. The margin in RB-Loss is default as 0.02.

##### C. FER with occlusion and variant pose in the wild

To address the occlusion and pose variant issues, we construct several test subsets with occlusion and pose annotations, i.e. Occlusion-FERPlus, Pose-FERPlus, Occlusion-AffectNet, Pose-AffectNet, Occlusion-RAF-DB, and Pose-RAF-DB. We

<sup>1</sup><http://dlib.net>

<sup>2</sup><https://pytorch.org/>

<sup>3</sup>[http://www.robots.ox.ac.uk/~vgg/software/vgg/\\$\\_face/](http://www.robots.ox.ac.uk/~vgg/software/vgg/$_face/)

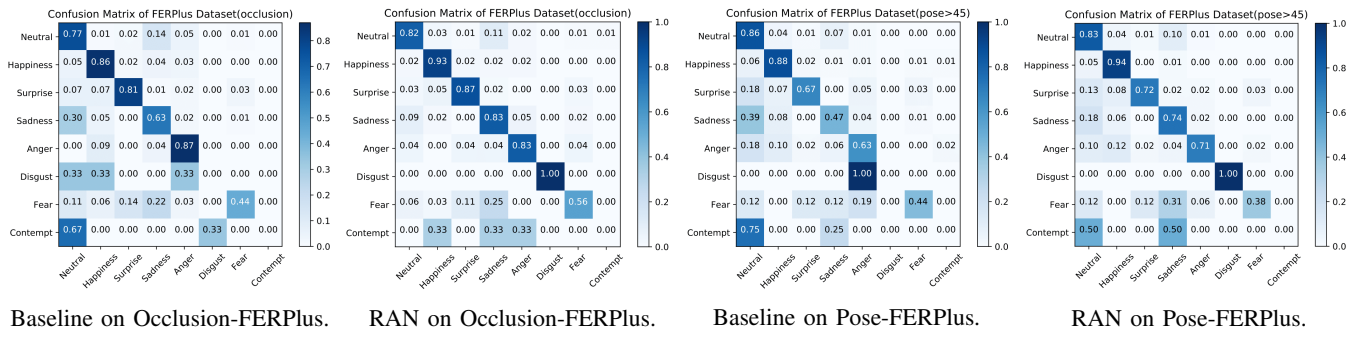


Fig. 4: The confusion matrices of baseline methods and our RAN on the Occlusion- and Pose-FERPlus test sets.

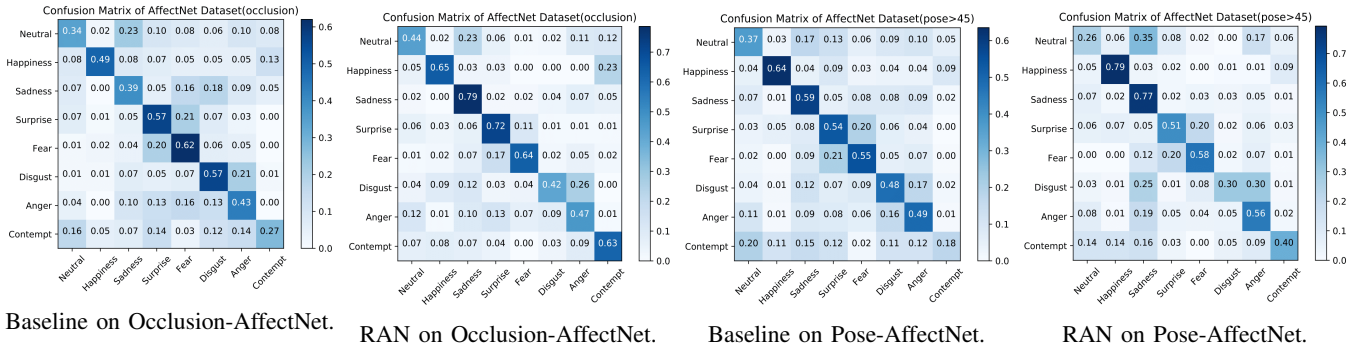


Fig. 5: The confusion matrices of baseline methods and our RAN on the Occlusion- and Pose-AffectNet test sets.

evaluate our RAN on the collected datasets with the default setting ( *i.e.* ResNet18 with alignment, RB-Loss and fixed cropping). We fine-tune the ResNet18 on original face images as baselines. Table II presents the comparison between the baselines and our method. Our RAN improves the baseline method significantly, with gains 10.3%, 10.02% and 2.53% on Occlusion-FERPlus, Occlusion-AffectNet and Occlusion-RAF-DB, respectively. On Pose-FERPlus, Pose-AffectNet and Pose-RAF-DB, the RAN also outperforms the baseline with a large margin. Specifically, with pose larger than 30 degrees, the gains are 4.12%, 3.09% and 2.70% on Pose-FERPlus, Pose-AffectNet and Pose-RAF-DB, respectively. The gains are improved to 4.9%, 5.4% and 2.05% with pose larger than 45 degrees. Overall, these results demonstrate the effectiveness of our proposed RAN on occlusion and variant pose FER data.

We present the confusion matrices of our RAN and these baselines in Figure 4 and Figure 5 to further investigate our improvements. We find that our RAN consistently boosts the “happiness”, “surprise”, and “sadness” categories on all the test sets. It may be explained by that these facial expressions have clear region features, such as action units of “Lip Corner Puller”, “Cheek Raiser”, and “Lip Corner Depressor”, which can be effectively captured by our RAN.

We also conduct a fair comparison on the recent occlusion test dataset: FED-RO [37]. We use the RAN with default setting, and train it using the same training data as [37]. We finally achieve **67.98%** which is clearly better than 66.5% of [37].

**Individual regions and their combination.** Since our RAN integrates several regions in a single network, we present the

TABLE III: The performance of individual regions with occlusion and variant pose conditions on FERPlus. ‘Aug. Training’ means that we augment the dataset by combining all the regions and original images and then train the network.

Region	Occlusion	Pose(30)	Pose(45)
Original ( $I_0$ )	73.33	78.11	75.50
$I_1$	67.43	74.27	71.40
$I_2$	64.13	72.22	70.30
$I_3$	72.22	78.48	76.84
$I_4$	72.8	78.54	77.00
$I_5$	74.54	78.63	75.35
Score Fusion ( $I_0 - I_5$ )	75.70	79.84	78.45
Aug. Training ( $I_0 - I_5$ )	79.92	81.24	79.26
<b>RAN (w RB-Loss)</b>	<b>83.63</b>	<b>82.23</b>	<b>80.40</b>

performance of individual regions and their score fusion on Occlusion- and Pose-FERPlus in Table III. To investigating if the improvement of our RAN only comes from augmented data, we also train a traditional model by mixing all the regions and the original images for data augmentation, which is called *aug. training*. We conclude that i) the performance of individual regions are comparable to each other except for the region  $I_1$  and  $I_2$ , ii) a naive score fusion (*i.e.* average) and mixing all the regions improve individual performance slightly, and iii) our RAN outperforms the score fusion and Aug. Training by a large margin. Compared to score fusion and fusion training, our RAN takes account of the importance of region features and also emphasises the most important region with RB-Loss.

**What is learned for occlusion and pose variant faces?**

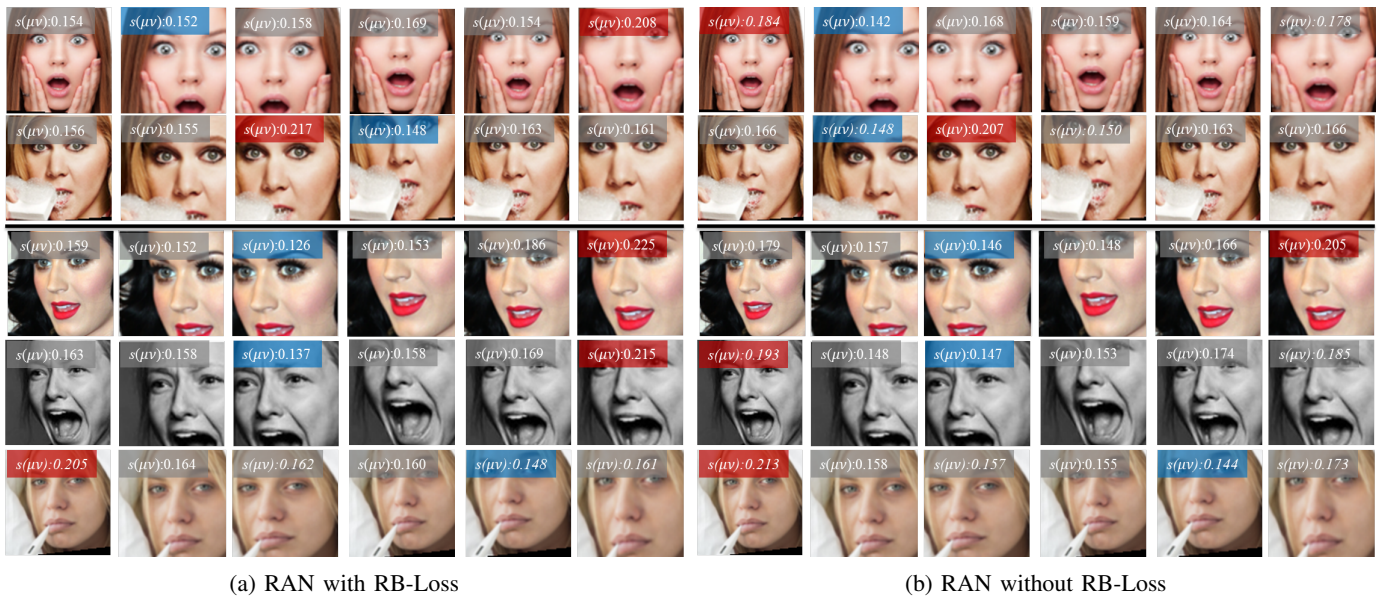


Fig. 6: Illustration of learned attention weights for different regions along with original faces.  $s(\cdot)$  denotes the softmax function. Red-filled boxes indicate the highest weights while blue-filled ones are the lowest weights. From left to right, the columns represent the original faces, regions  $I_1$  to  $I_5$ . Note that the left and right figures show the weights obtained with and without the RB-Loss respectively. Better viewed in PDF.

To better explore our RAN, we illustrate the final attention weights for several examples with RB-Loss and without RB-Loss in Figure 6(a) and Figure 6(b), respectively. Occlusion examples are shown in the first two rows, and pose examples in the third and fourth rows. We also show a bad example in the last row.

For the occlusion examples, our RAN with RB-Loss gets the highest weight on the small center crop (i.e.  $I_5$ ) for the first example. It makes sense since this image suffers from the left and right occlusion. In the second example which suffers from bottom-left occlusion, the RAN with RB-Loss automatically assigns the highest weight to the up-right region while suppresses the bottom-left region. For both pose-variant examples in the third and fourth rows, our RAN with RB-Loss gets high attention weights on center regions while gets low weights on the up-right regions. This may be explained by that the up-right regions contain the most of irrelevant information on the near-profile faces. With RB-Loss and RAN, the original faces get relatively average attention weights among all the examples. For a bad example in the last row, our RAN with RB-Loss does not assign the regions with high weights. It may illustrate that the tiny occlusion can not impact the FER system a lot. The attention weights of tiny occlusion face are more likely to be random. Compared to the RB-Loss case, though RAN without RB-Loss can also assign different attention weights similarly, the weights for all the regions from RAN without RB-Loss are smoother. In addition, the original image prefers to have the highest weight without RB-Loss.

#### D. Ablation study on FERPlus and AffectNet

To validate the generality of our method, we conduct an ablation study on the full test set of FERPlus and the full

TABLE IV: Evaluation of all components of our RAN along with face alignment on FERPlus.

Align	Self-att.	Relation-att.	RB-Loss	Accuracy
				86.50
	✓			86.90
	✓	✓		87.63
	✓	✓	✓	<b>87.85</b>
✓				87.60
✓	✓			87.80
✓	✓	✓		88.23
✓	✓	✓	✓	<b>88.55</b>

validation set of AffectNet with default setting. Face alignment is a standard pre-processing method for face analysis, while a few works do not utilize [5], [48] for FER task. Here we also study the effect of face alignment.

**Attention modules.** We first study the attention modules of our RAN without using RB-Loss. The evaluation results on FERPlus and AffectNet are presented on Table IV and Table V, respectively. On FERPlus without face alignment, the self-attention ( $F_m$  in Eq. (3)) improves the baseline by 0.4%. Adding the relation-attention module, our method outperforms the baseline by 1.13% and 3.05% on FERPlus and AffectNet without face alignment. Face alignment is found to significantly boost the baseline method on both datasets, while its effect is limited when using our proposed RAN. This can be explained by that our method implicitly learns to align facial regions with the attention mechanism as that in machine translation [4]. With face alignment, our attention modules improve the baselines by 0.83% and 1.85% on FERPlus and AffectNet, respectively.

**Region biased loss.** The RB-Loss is added to the self-attention module with margin 0.02 by default. From Table IV



TABLE V: Evaluation of all components of our RAN along with face alignment on AffectNet **without oversampling**.

Align	Self-att.	Relation-att.	RB-Loss	Accuracy
				49.00
	✓	✓		52.05
✓	✓	✓	✓	<b>52.97</b>
✓	✓	✓		50.32
✓	✓	✓		52.17
✓	✓	✓	✓	52.50

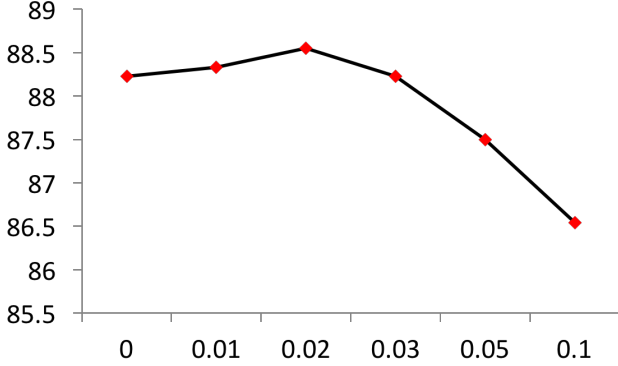


Fig. 7: The evaluaiton of the margin ( $\alpha$ ) in RB-Loss on FERPlus.

and Table V, we can see that the designed RB-Loss further improves performance on both FERPlus and AffectNet consistently. Specifically, the improvement on AffectNet without face alignment is 0.92%. With oversampling, our RAN with RB-Loss achieves **59.5%** on the validation set of AffectNet. It is worth noting that RB-Loss does not increase computational cost in testing.

We also evaluate the parameter  $\alpha$  of RB-Loss in Figure 7. Increasing  $\alpha$  from 0 to 0.02 gradually improves the performance while larger  $\alpha$  leads to fast degradation, which indicates the original image is also important for FER. As the mater of fact, the result of this experiment is part of our motivation to keep the original face image for our method.

**Evaluation of individual regions and different fusion schemes.** We conduct an evaluation of individual regions and different fusion schemes on the full FERPlus test datasets without face alignment. For the fusion schemes, we mainly consider three popular methods, namely feature concatenation, feature average pooling, and score fusion (i.e. score average). The evaluation results are shown in Figure 8. Several observations can be concluded as follows. First, all the individual crops are inferior to the original image which indicates the performance gain is not from special enlarged crops. Second, compared to the original image, there is no obvious improvement by concatenating and averaging region features. Third, score fusion slightly improves the baseline by 0.54% while our RAN outperforms the baseline by 1.35%.

**Evaluation of region generation strategies.** We evaluate the fixed cropping, landmark-based cropping, and random cropping methods on FERPlus with the default setting for other parameters, the results are shown in Figure 9. For random cropping, we randomly generate 3 regions for each

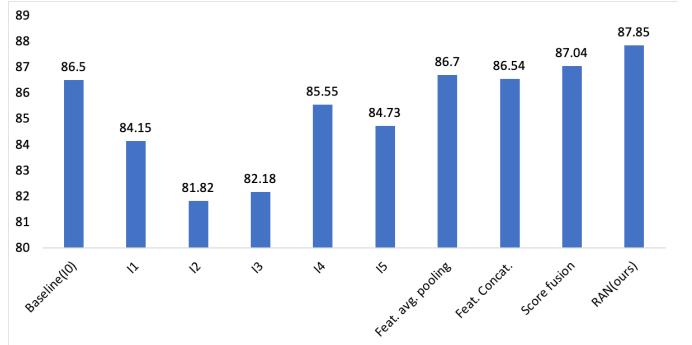


Fig. 8: Performce comparison of individual regions and different aggregation schemes on FERPlus.

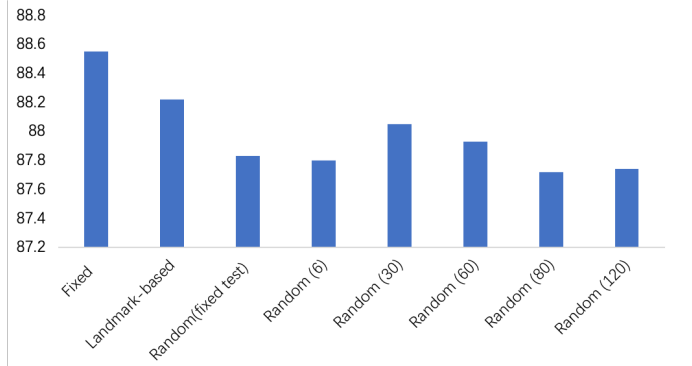


Fig. 9: Evaluation of different region generation strategies and the number of random regions.

image in each training iteration while generate 6, 30, 60, 80, 120 random regions for test evaluation several times. For landmark-based cropping, we set the radius as 0.4 of the side of image which ensures a similar size as fixed cropping. The fixed cropping strategy consistently outperforms the random cropping even dozens of times more regions are used. The landmark-based cropping performs slightly worse than the fixed cropping. Training with random cropping yet testing with the same fixed cropping has limited effect for random region cropping. Increasing the crops boosts performance in the beginning while degrades after 30 crops. This may be explained by that increasing crops leads to too many sub-optimized regions and they dominate the final representation.

**Evaluation of region size.** To explore the impact of region sizes for our RAN, we evaluate the region size of fixed cropping scheme on FERPlus with other parameter as default. Since five regions with different sizes are cropped in our default setting, we evaluate these region sizes using a ratio from low to high compared to the default sizes. The evaluation results are shown in Figure 10. The performance degrades significantly with the ratio reducing to 0.4. Increasing (i.e. ratio:1.1) size upon the default one slightly reduces the performance. It may be explained that the regions of  $I_4$  and  $I_5$  almost degrade to the original image, and the information gain from enlarging regions disappeared if too large regions are used.

**Evaluation of inference time.** Since our RAN has five

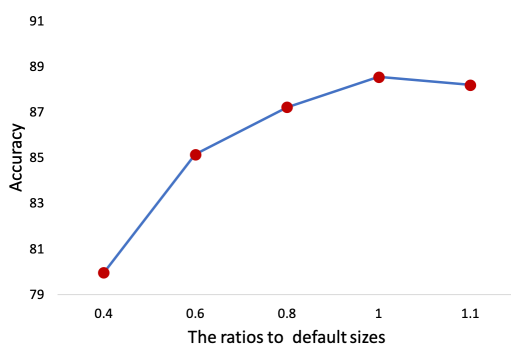


Fig. 10: Evaluation of region sizes on FERPlus. The ratios are compared to the default setting.

TABLE VI: Comparison to the state-of-the-art results on the FERPlus dataset.\*These results are trained using label distribution as supervision.

Method	Network	Pre-trained Dataset	Year	Performance
[5]*	VGG13	/	2016	85.1
[27]	ResNet18+VGG16	/	2017	87.4
[1]*	SeNet50	VGG_Face2	2018	88.8
RAN-ResNet18	ResNet18	MS_Celeb_1M	2019	88.55
RAN-VGG16*	VGG16	VGG_Face	2019	<b>89.16</b>

times feedforward operations than the baseline, we investigate the inference time on FERPlus test set. We evaluate the average per-image inference time and the run-time is obtained on a TITAN 1080ti GPU of Linux Cluster with a 2.6 GHz Intel(R) Xeon(R) E5-2690 CPU. The average inference time of RAN and the baseline are 0.025s and 0.006s, respectively. Due to the powerful parallel ability of GPU, the increasing time is not linear to the number of regions.

#### E. Comparison with the state-of-the-art methods

In this section, we compare our best results to several state-of-the-art methods on FERPlus, AffectNet, SFEW, and RAF-DB.

**Comparison on FERPlus.** We compare our RAN to several state-of-the-art methods on the FERPlus dataset in Table VI. Both [5] and [1] leverage the label distribution for each face as supervision. [1] pretrains a SeNet50 [26] on VGGFace2.0 [10] which includes amount of large-pose faces. With the KLDiv loss and label distribution supervision, we fine-tune the public VGGFace model (VGG16 pretrained on VGGFace1.0) with our RAN and achieve 89.16% which is a new state of the art to our knowledge.

**Comparison on AffectNet.** Table VII presents the comparison on AffectNet. We obtain 52.97% and 59.5% without and with oversampling, respectively. It is worth noting that [48] only achieves 47% with upsampling and [66] uses one more large-scale FER dataset and 80 layers ResNet for training with elaborated loss weights on them.

**Comparison on SFEW.** Table VIII presents the comparison on SFEW. [9] applies a small CNN with an island loss which is the combination of the Center loss [60] and an inter-class loss. [65] ensembles multiple CNNs with each CNN model initialized randomly or pretrained on FER2013. Our RAN with

TABLE VII: Comparison to the state-of-the-art results on the AffectNet dataset.<sup>+</sup>Oversampling is used for a final performance report since AffectNet is imbalanced. <sup>‡</sup>RAF-DB is added into training data.

Method	Network	Year	Performance
Up-Sampling [48]	AlexNet	2018	47.0
Weighted-Loss [48]	AlexNet	2018	58.0
[66] <sup>‡</sup>	ResNet80	2018	55.71
RAN-ResNet18	ResNet18	2019	52.97
RAN-ResNet18 <sup>+</sup>	ResNet18	2019	<b>59.5</b>

TABLE VIII: Comparison to the state-of-the-art results on the SFEW dataset.

Method	Pre-trained Dataset	Year	Performance
Island Loss [9]	FER2013	2018	52.52
Identity-aware CNN [46]	FER2013	2017	50.98
Multiple deep CNNs [65]	FER2013	2015	55.96
RAN-ResNet18	MS_Celeb_1M	2019	54.19
RAN(VGG16+ResNet18)	MS_Celeb_1M	2019	<b>56.4</b>

single model achieves 54.19% on the validation set which is the best single model to our best of knowledge. Since model ensemble is popular on SFEW, we also conduct a naive model fusion by averaging the scores of ResNet18 and VGG16 which obtains 56.4%.

**Comparison on RAF-DB.** Table IX presents the comparison on RAF-DB. RAF-DB is a latest facial expression dataset which not only has basic emotion categories but also compound categories. We report the overall accuracy on the basic emotion categories. [34] introduces the RAF-DB dataset and uses a locality-preserving loss for network training. [37] leverages patch-based attention networks and global networks. Our proposed RAN achieves **86.9%** on RAF-DB with default setting, which are 2.77% and 1.83% better than DLP-CNN [34] and [37], respectively.

## V. CONCLUSION

In this paper, we address the facial expression recognition in the real-world occlusion and pose-variant conditions. We build several new FER test datasets on these conditions, and propose the Region Attention Network (RAN) which adaptively adjusts the importance of facial parts. We further design a region Biased loss (RB-Loss) function to encourage high attention weight for the most important region. We evaluate our method on the collected datasets and make extensive studies on FER-Plus and AffectNet. Our proposed method achieves state-of-the-art results on FERPlus, SFEW, RAF-DB, and AffectNet.

## REFERENCES

- [1] S. Albanie, A. Nagrani, A. Vedaldi, and A. Zisserman. Emotion recognition in speech using cross-modal transfer in the wild. In *ACM Multimedia*, 2018. 1, 10
- [2] Samuel Albanie, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Emotion recognition in speech using cross-modal transfer in the wild. *arXiv preprint arXiv:1808.05561*, 2018. 2
- [3] Brandon Amos, Bartosz Ludwiczuk, and Mahadev Satyanarayanan. Openface: A general-purpose face recognition library with mobile applications. Technical report, CMU-CS-16-118, CMU School of Computer Science, 2016. 1, 2, 5

TABLE IX: Comparison to the state-of-the-art results on the RAF-DB dataset.

Method	Network	Year	Performance
DLP-CNN [34]	8-layer baseDCNN	2019	84.13
gACNN [37]	VGG16	2018	85.07
RAN-ResNet18	ResNet18	2019	<b>86.90</b>

- [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014. **2, 8**
- [5] Emad Barsoum, Cha Zhang, Cristian Canton Ferrer, and Zhengyou Zhang. Training deep networks for facial expression recognition with crowd-sourced label distribution. In *International Conference on Multimodal Interaction*, pages 279–283, 2016. **1, 5, 6, 8, 10**
- [6] Júlio César Batista, Vítor Albiero, Olga RP Bellon, and Luciano Silva. Aumpnet: simultaneous action units detection and intensity estimation on multipose facial images using a single convolutional neural network. In *FG*, pages 866–871, 2017. **1**
- [7] Jerry D Boucher and Paul Ekman. Facial areas and emotional information. *Journal of communication*, 25(2):21–29, 1975. **1, 2, 4**
- [8] Fabrice Bourel, Claude C Chibellushi, and Adrian A Low. Recognition of facial expressions in the presence of occlusion. In *BMVC*, pages 1–10, 2001. **1**
- [9] Jie Cai, Zibo Meng, Ahmed Shehab Khan, Zhiyuan Li, James O'Reilly, and Yan Tong. Island loss for learning discriminative features in facial expression recognition. In *FG*, pages 302–309. IEEE, 2018. **10**
- [10] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 67–74. IEEE, 2018. **10**
- [11] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory networks for machine reading. *CoRR*, abs/1601.06733, 2016. **2**
- [12] Shane F Cotter. Sparse representation for accurate classification of corrupted and occluded facial expressions. In *ICASSP*, pages 838–841, 2010. **1, 2**
- [13] Shane F Cotter. Weighted voting of sparse representation classifiers for facial expression recognition. In *Signal Processing European Conference*, pages 1164–1168, 2010. **1, 2**
- [14] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005. **2**
- [15] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon. Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark. In *IJCV Workshops*, 2011. **6**
- [16] Abhinav Dhall, Roland Goecke, Simon Lucey, Tom Gedeon, et al. Collecting large, richly annotated facial-expression databases from movies. *IEEE Multimedia*, 19(3):34–41, 2012. **6**
- [17] Abhinav Dhall, Amanjot Kaur, Roland Goecke, and Tom Gedeon. EmotiW 2018: Audio-video, student engagement and group-level affect prediction. In *International Conference on Multimodal Interaction*, pages 653–656. ACM, 2018. **1**
- [18] Abhinav Dhall, OV Ramana Murthy, Roland Goecke, Jyoti Joshi, and Tom Gedeon. Video and image based emotion recognition challenges in the wild: EmotiW 2015. In *International Conference on Multimodal Interaction*, pages 423–426. ACM, 2015. **6**
- [19] Changxing Ding and Dacheng Tao. A comprehensive survey on pose-invariant face recognition. *ACM Transactions on intelligent systems and technology (TIST)*, 7(3):37, 2016. **1**
- [20] H. Ding, S. K. Zhou, and R. Chellappa. Facenet2expnet: Regularizing a deep face recognition net for expression recognition. In *2017 12th IEEE International Conference on Automatic Face Gesture Recognition (FG 2017)*, pages 118–126, May 2017. **2**
- [21] C Fabian Benitez-Quiroz, Ramprakash Srinivasan, and Aleix M Martinez. Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In *CVPR*, pages 5562–5570, 2016. **1**
- [22] B. Fasel. Robust face analysis using convolutional neural networks. In *ICPR*, pages 40–43, 2002. **2**
- [23] Ian J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hammer, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. Challenges in representation learning: A report on three machine learning contests. In *International Conference on Neural Information Processing*, pages 117–124. Springer, 2013. **6**
- [24] SL Happy and Aurobinda Routray. Automatic facial expression recognition using features of salient facial patches. *Transactions on Affective Computing*, 6(1):1–12, 2015. **1**
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. **6**
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. **10**
- [27] Christina Huang. Combining convolutional neural networks for emotion recognition. In *URTC, 2017 IEEE MIT*, pages 1–4. IEEE, 2017. **10**
- [28] Heechul Jung, Sihaeng Lee, Junho Yim, Sunjeong Park, and Junmo Kim. Joint fine-tuning in deep neural networks for facial expression recognition. In *ICCV*, pages 2983–2991, 2015. **1**
- [29] Samira Ebrahimi Kahou, Christopher Pal, Xavier Bouthillier, Pierre Froumenty, Roland Memisevic, Pascal Vincent, Aaron Courville, Yoshua Bengio, and Raul Chandias Ferrari. Combining modality specific deep neural networks for emotion recognition in video. In *International Conference on Multimodal Interaction*, pages 543–550, 2013. **1, 2**
- [30] Takeo Kanade, Yingli Tian, and Jeffrey F Cohn. Comprehensive database for facial expression analysis. In *FG*, page 46. IEEE, 2000. **1**
- [31] Vahid Kazemi and Josephine Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, pages 1867–1874, 2014. **6**
- [32] Irene Kotsia, Ioan Buciu, and Ioannis Pitas. An analysis of facial expression recognition under partial facial image occlusion. *Image and Vision Computing*, 26(7):1052–1067, 2008. **1, 2**
- [33] Gil Levi and Tal Hassner. Emotion recognition in the wild via convolutional neural networks and mapped binary patterns. In *Proceedings of the 2015 ACM on international conference on multimodal interaction*, pages 503–510. ACM, 2015. **2**
- [34] S. Li and W. Deng. Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. *IEEE Transactions on Image Processing*, 28(1):356–370, Jan 2019. **5, 6, 10, 11**
- [35] Shan Li and Weihong Deng. Deep facial expression recognition: A survey. *CoRR*, abs/1804.08348, 2018. **1**
- [36] Shan Li, Weihong Deng, and JunPing Du. Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild. In *CVPR*, pages 2584–2593, 2017. **1**
- [37] Y. Li, J. Zeng, S. Shan, and X. Chen. Occlusion aware facial expression recognition using cnn with attention mechanism. *IEEE Transactions on Image Processing*, 28(5):2439–2450, May 2019. **2, 5, 7, 10, 11**
- [38] Shengcai Liao, Anil K Jain, and Stan Z Li. Partial face recognition: Alignment-free approach. *TPAMI*, 35(5):1193–1205, 2013. **1**
- [39] Chengjun Liu and H. Wechsler. Gabor feature based classification using the enhanced fisher linear discriminant model for face recognition. *IEEE Transactions on Image Processing*, 11(4):467–476, April 2002. **2**
- [40] Mengyi Liu, Shaolin Li, Shiguang Shan, and Xilin Chen. Au-inspired deep networks for facial expression feature learning. *Neurocomputing*, 159(C):126–136, 2015. **1, 2**
- [41] S. Liu, Y. Zhang, K. Liu, and Y. Li. Facial expression recognition under partial occlusion based on gabor multi-orientation features fusion and local gabor binary pattern histogram sequence. In *International Conference on Intelligent Information Hiding and Multimedia Signal Processing*, 2013. **2**
- [42] Xiang Long, Chuang Gan, Gerard de Melo, Jiajun Wu, Xiao Liu, and Shilei Wen. Attention clusters: Purely attention based local feature integration for video classification. *CoRR*, abs/1711.09550, 2017. **3**
- [43] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015. **4**
- [44] Michael J Lyons, Shigeru Akamatsu, Miyuki Kamachi, Jiro Gyoba, and Julien Budynek. The japanese female facial expression (jaffe) database. In *FG*, pages 14–16, 1998. **1**
- [45] Debin Meng, Xiaojiang Peng, Kai Wang, and Yu Qiao. frame attention networks for facial expression recognition in videos. *arXiv preprint arXiv:1907.00193*, 2019. **2**
- [46] Zibo Meng, Ping Liu, Jie Cai, Shizhong Han, and Yan Tong. Identity-aware convolutional neural network for facial expression recognition. In *FG*, pages 558–565. IEEE, 2017. **10**
- [47] Volodymyr Mnih, Nicolas Heess, Alex Graves, and koray kavukcuoglu. Recurrent models of visual attention. In *NIPS*, pages 2204–2212. 2014. **2**
- [48] Ali Mollahosseini, Behzad Hasani, and Mohammad H Mahoor. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *Transactions on Affective Computing*, 2017. **1, 5, 6, 8, 10**

- [49] Pauline C. Ng and Steven Henikoff. Sift: predicting amino acid changes that affect protein function. *Nucleic Acids Research*, 31(13):3812–3814, 2003. 2
- [50] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015. 6
- [51] O. Rudovic, M. Pantic, and I. Patras. Coupled gaussian processes for pose-invariant facial expression recognition. *TPAMI*, 35(6):1357–1369, June 2013. 2
- [52] Caifeng Shan, Shaogang Gong, and Peter W. McOwan. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing*, 27(6):803 – 816, 2009. 2
- [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [54] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000 classes. In *CVPR*, pages 1891–1898, 2014. 5
- [55] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 2
- [56] Lianzhi Tan, Kaipeng Zhang, Kai Wang, Xiaoxing Zeng, Xiaojiang Peng, and Yu Qiao. Group emotion recognition with individual facial emotion cnns and global image based cnns. In *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, pages 549–552. ACM, 2017. 2
- [57] Yichuan Tang. Deep learning using linear support vector machines. *Computer Science*, 2013. 1, 2
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008. Curran Associates, Inc., 2017. 2
- [59] Jianfeng Wang, Ye Yuan, and Gang Yu. Face attention network: An effective face detector for the occluded faces. *CoRR*, abs/1711.07246, 2017. 3
- [60] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016. 10
- [61] Flood Sung Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. 4
- [62] Jiaolong Yang, Peiran Ren, Dong Chen, Fang Wen, Hongdong Li, and Gang Hua. Neural aggregation network for video face recognition. *CoRR*, abs/1603.05474, 2016. 3
- [63] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 2
- [64] Galit Yovel and Brad Duchaine. Specialized face perception mechanisms extract both part and spacing information: Evidence from developmental prosopagnosia. *Journal of Cognitive Neuroscience*, 18(4):580–593, 2006. PMID: 16768361. 1
- [65] Zhiding Yu and Cha Zhang. Image based static facial expression recognition with multiple deep network learning. In *International Conference on Multimodal Interaction*, pages 435–442. ACM, 2015. 10
- [66] Jiabei Zeng, Shiguang Shan, and Xilin Chen. Facial expression recognition with inconsistently annotated datasets. In *ECCV*, pages 222–37, 2018. 10
- [67] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, Oct 2016. 2
- [68] Kaipeng Zhang, Lianzhi Tan, Zhifeng Li, and Yu Qiao. Gender and smile classification using deep convolutional neural networks. In *CVPR Workshops*, pages 34–38, 2016. 5
- [69] Ligang Zhang, Dian Tjondronegoro, and Vinod Chandran. Facial expression recognition experiments with data from television broadcasts and the world wide web. *Image and Vision Computing*, 32(2):107–119, 2014. 1
- [70] Ligang Zhang, Dian Tjondronegoro, and Vinod Chandran. Random gabor based templates for facial expression recognition in images with facial occlusion. *Neurocomputing*, 145:451–464, 2014. 1
- [71] Shiqing Zhang, Xiaoming Zhao, and Bicheng Lei. Robust facial expression recognition via compressive sensing. *Sensors*, 12(3):3747–3761, 2012. 1
- [72] Xiangyun Zhao, Xiaodan Liang, Luoqi Liu, Teng Li, Yugang Han, Nuno Vasconcelos, and Shuicheng Yan. Peak-piloted deep network for facial expression recognition. In *European conference on computer vision*, pages 425–442. Springer, 2016. 2