

# Boosting Facial Expression Recognition by A Semi-Supervised Progressive Teacher

Jing Jiang and Weihong Deng, *Member, IEEE*

**Abstract**—In this paper, we aim to improve the performance of in-the-wild Facial Expression Recognition (FER) by exploiting semi-supervised learning. Large-scale labeled data and deep learning methods have greatly improved the performance of image recognition. However, the performance of FER is still not ideal due to the lack of training data and incorrect annotations (e.g., label noises). Among existing in-the-wild FER datasets, reliable ones contain insufficient data to train robust deep models while large-scale ones are annotated in lower quality. To address this problem, we propose a semi-supervised learning algorithm named Progressive Teacher (PT) to utilize reliable FER datasets as well as large-scale unlabeled expression images for effective training. On the one hand, PT introduces semi-supervised learning method to relieve the shortage of data in FER. On the other hand, it selects useful labeled training samples automatically and progressively to alleviate label noise. PT uses selected clean labeled data for computing the supervised classification loss and unlabeled data for unsupervised consistency loss. Experiments on widely-used databases RAF-DB and FERPlus validate the effectiveness of our method, which achieves state-of-the-art performance with accuracy of 89.57% on RAF-DB. Additionally, when the synthetic noise rate reaches even 30%, the performance of our PT algorithm only degrades by 4.37%.

**Index Terms**—Facial Expression Recognition, Semi-supervised Learning, Label noise.



## 1 INTRODUCTION

Facial expression is one of the most natural and universal way for human beings to convey their emotional states and intentions. Facial expression recognition (FER) has become a research hotspot in the field of computer vision since it's significant for human-computer interaction applications, such as remote education and driving fatigue monitoring devices. Researches on FER so far have mostly focused on discrete basic expression categories, i.e., anger, disgust, fearful, happy, sad and surprised. Recently, compound expressions have also been considered to describe more fine-grained emotions. According to data sources, databases for FER can be divided into two groups, which are lab-controlled and in-the-wild ones respectively. Lab-controlled expression datasets are conducted under the laboratory environment that subjects are asked to make specific expressions. This leads to accurate annotations but insufficient data. In-the-wild datasets consist of images collected from the Internet and manual annotations, so they usually have larger data quantity than lab-controlled ones due to easier data collection. However, it's time-consuming and resource-consuming to annotate a large-scale reliable dataset, which causes the contradiction between data quality and quantity.

Early studies use lab-controlled datasets and hand-crafted features such as local binary patterns (LBP), histogram of gradient (HoG), and scale invariant feature transform (SIFT) to extract discriminative information and then train classifiers. This seems to work well because of the simple patterns of data. But the classifier will have poor perfor-

mance on unseen data, especially images in real scenarios. Afterwards, significant progress has been made towards improving the performance of FER, especially researches based on deep learning method [1], [2], [3], [4]. With the development of deep learning, Convolutional Neural Networks (CNNs) remarkably benefit image recognition task. Large amount of high-quality data is crucial for training a robust CNN model. Considering training data, in-the-wild datasets [5], [6], [7] are more suitable for FER in real-world condition than lab-controlled ones [8], [9], [10]. On the one hand, in-the-wild datasets usually have more training samples. On the other hand, more complex expression diversity makes the model more robust and adaptive. Even so, real-world FER still faces some challenges. Firstly, in-the-wild datasets may still have insufficient samples to train a robust deep neural network. Secondly, it's easy to collect expression-related images from Internet but it's hard to annotate them accurately. So the datasets may have inconsistent labels or incorrect labels (noisy labels) due to the uncertainty of samples [4]. These two issues cause over-fitting problem and thus limit the recognition performance, especially for data-driven deep learning based FER.

Considering existing expression databases, we aim to fully utilize the reliable in-the-wild dataset as well as exploiting the effect of large-scale dataset with poor annotations. For example, AffectNet [6] is by far the largest database for FER which contains more than one millions of expression-related images. Among them about 280,000 training images are annotated into seven basic expressions (neutral is also considered), but each image is labeled by only one human-encoder so that the quality of annotation is not ensured. Therefore, the deep model trained with this dataset may have poor generalization performance. We show some apparently mislabeled samples in AffectNet in Fig.1. Reliable in-the-wild FER datasets such as RAF-DB [5]

- *The authors are with the Pattern Recognition and Intelligent System Laboratory, School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, 100876, China. Weihong Deng is the corresponding author. E-mail: {jiangjing1998, whdeng}@bupt.edu.cn.*
- *This work was supported by the National Natural Science Foundation of China under Grants No. 61871052*

and FERPlus [7] contains much less data which also limits the strength of deep CNN. RAF-DB contains 12,271 training samples and each of them is labeled by 40 independent individuals. FERPlus consists of more than 20,000 training data and each of them is annotated by 10 experts. Label noises are relatively rare in these two datasets. We list the data distribution of above three datasets in Table1.

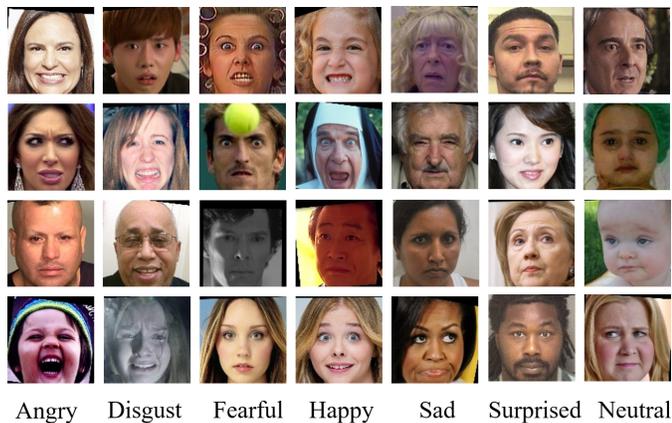


Fig. 1. Examples of not accurately annotated images in AffectNet.

TABLE 1  
The data distribution of RAF-DB, FERPlus and AffectNet among seven expressions.

	Ang	Dis	Fea	Hap	Neu	Sad	Sur
RAF-DB	705	717	281	4772	2524	1982	1290
FERPlus	2399	175	648	7410	9365	3403	3378
AffectNet	24882	3803	6378	134415	74874	25459	14090

Since it's difficult and time-consuming to annotate a large-scale in-the-wild FER dataset accurately, we can utilize the large amount of face images without annotations to help train a more robust model. Semi-supervised learning points out that the classifier which uses auxiliary unlabeled data can outperform that only uses small-size labeled ones. Therefore, we intend to take maximum advantage of the large-scale expression images in AffectNet as unlabeled data to boost the recognition performance.

To this end, we propose a semi-supervised deep learning framework named Progressive Teacher (PT) to address the shortage of data and label noise simultaneously. Progressive Teacher adopts semi-supervised learning to relieve the lack of training data in FER that auxiliary large amount of unlabeled data are utilized to boost the performance. Inspired by semi-supervised algorithm Mean Teacher [11], we follow the architecture of teacher-student model and make the network work in semi-supervised manner. Different from previous works on FER [12], [13] that unlabeled data are used to pre-train the network or annotated automatically to enlarge training set, we use both labeled and unlabeled data for computing the overall loss in a unified framework. Fig.2 shows how the large-scale unlabeled data are utilized. The teacher model has stronger learning ability and guides the training of student model. Student model improves by

learning to produce consistent outputs with the teacher model besides optimizing the classification loss of labeled data. Considering it costs much to construct a reliable in-the-wild FER database, we expect the network to be robust to label noises so that it can work even in noisy datasets. However, when the labeled dataset contains certain noisy labels, the negative effect of overfitting will gradually accumulate in this teacher-student framework such as Mean Teacher. Our Progressive Teacher alleviate this phenomenon by selecting useful labeled training samples automatically and progressively to tackle label noises. As illustrated in Fig.2, in each iteration, the teacher model will filter out a portion of labeled samples with increasing rate which are believed to be noisy ones, so as to prevent the student model from learning from their inaccurate annotations. Those samples are thereby regarded as unlabeled and used for consistency loss. In this case, labeled data are effectively utilized to prevent from overfitting and unlabeled data also improve the generalization performance in the meanwhile.

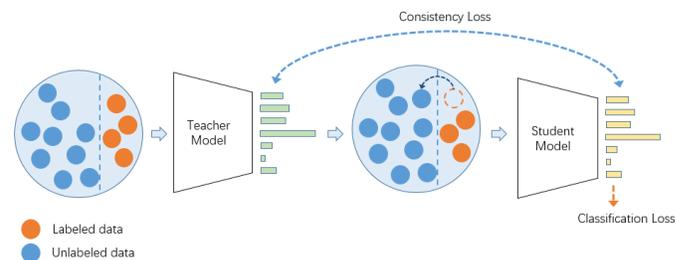


Fig. 2. An illustration of Progressive Teacher.

Overall, our contributions can be summarized as follows,

1. We analyze existing in-the-wild FER databases and aim to tackle the problem that reliable datasets contain insufficient labeled samples to train a robust deep neural network. Therefore, a semi-supervised algorithm Progressive Teacher (PT) is proposed. PT utilizes auxiliary large-scale unlabeled expression images for effective training in a semi-supervised manner to boost the recognition accuracy. Different from Mean Teacher, the teacher not only provides better logits and guides student to produce similar outputs, but also selects potential clean samples for student models to learn. Additionally, We use two pairs of teacher-student models and adopt a cross-guidance mechanism.
2. To alleviate the performance drop caused by noisy labels in traditional semi-supervised methods such as Mean Teacher, the proposed PT prevents from overfitting noisy samples by selecting clean labeled samples automatically and progressively.
3. We extensively validate our PT algorithm on in-the-wild datasets RAF-DB and FERPlus. The results indicate that our semi-supervised method improves the recognition performance. Remarkably, we achieve the highest accuracy of 89.57% on RAF-DB as far as we know. We also evaluate the effect of PT to tackle label noises on synthetic noisy datasets. It alleviates the performance drop evidently. When noisy samples get more, the effect of PT gets more remarkable. The experiment on real-world noisy dataset AffectNet also indicates the effectiveness of PT.

## 2 RELATED WORK

In this section, we mainly discuss methods that are related to facial expression recognition, semi-supervised learning and how to tackle with label noise and face uncertainties.

### 2.1 Facial expression recognition

As pointed out in [14], deep learning based facial expression recognition can be performed by three major steps: data pre-processing, deep feature learning and deep feature classification. We get aligned face regions from original images in the pre-processing stage, and then feed the aligned images into a deep neural network, for example, Convolutional Neural Network(CNN), to obtain discriminative features and then classification will be conducted almost simultaneously. Deep learning based FER requires large amount of training samples to avoid the over-fitting problem. However, reliable in-the-wild FER databases, such as RAF-DB [5] and FERPlus [7] usually have insufficient data to learn discriminative deep features. The annotation quality of the large-scale database AffectNet [6] is not ensured due to large data quantity, which limits the strength of CNN.

To tackle the shortage of training data in FER, some works [15], [16] utilize face recognition datasets to pre-train the network and then fine-tune it on the expression dataset. To eliminate the effect of face-dominated information reserved in the pre-trained face network, FaceNet2ExpNet [1] was proposed to make the pre-trained model provide a good initialization. Considering utilizing more expression datasets to enlarge the amount of data, IPA2LT [2] has found that the performance of FER can't be improved and even degrades by merging multiple datasets directly due to inconsistent annotation. To avoid this effect, IPA2LT uses multiple models trained on different datasets to discover the latent truth of input images. Recently, a omni-supervised FER baseline [13] was proposed to make use of auxiliary large-scale unlabeled images. It first trains a primitive model using a small number of labeled samples, then use it to select high confident unlabeled samples by feature-based similarity comparison. The enlarged dataset is proved to boost the recognition performance. It also adopts a dataset distillation strategy to distill and compress the useful knowledge from the selected auxiliary samples in order to reduce computation resources. [17] proposed a weakly supervised learning technique, which uses unlabeled data with high confidence scores as labeled ones and train the network in a supervised way. Similar to [13], we aim to exploit the effect of large-scale unlabeled images, but in an intrinsically semi-supervised manner. Our method is an extension of teacher-student model. [18] also employ the teacher-student training strategy and it treats model trained on fully-visible faces as teacher and model trained on occluded faces as student to learn discriminative embeddings for classifying expressions under occlusion. Knowledge distillation and triplet loss are utilized for learning embeddings.

### 2.2 Semi-supervised Learning

The success of deep learning bases on large amount of well-annotated data. However, constructing a large-scale dataset with high-quality annotation is time-consuming and labor-expensive. Semi-supervised learning tackles the shortage

of labeled training data and improves learning ability. As illustrated in [19], semi-supervised learning uses labelled as well as unlabeled data to perform certain learning tasks such as classification. Some methods [20], [21] which can be classes as self-training use base classifiers to obtain predictions for unlabeled data and these pseudo-labelled data are therefore used for supervised training. In [20], unlabeled images are labeled as the most confident of their predictions. The unlabeled data are weighted to fit the training process since pseudo labels are not reliable enough in the beginning. [21] balances the influence of unlabeled data by giving pseudo label when the maximum prediction confidence is greater than a threshold. Co-training [22] and Tri-training [23] are extensions of self-training which learn from multiple supervised classifiers. Some methods [11], [24], [25] based on deep learning make the network work in semi-supervised fashion instead of giving the unlabeled data pseudo labels in advance.  $\Pi$ -model [25] exploits unlabeled data by constraining the predictions of two models to be consistent under different data augmentation and dropout. The total loss consists of two parts, one is the classification loss of labeled data, the other is the consistency loss of all data without labels. Temporal ensembling [25] is an extension of  $\Pi$ -model, it penalizes the difference in the network outputs at different points in time during the training process. It's an implicit teacher-student model which constructs a better target by ensembling the outputs during training to regularize the student model. Compared to [25], Mean teacher [11] is an explicit teacher-student model which constructs a better teacher model by averaging the weights at each training iteration. Inspired by [11], we also make a good teacher model by weight-averaging and regularize the student model by penalizing the difference in network outputs with teacher model. But we use two pairs of teacher-student models, assisted by a cross-guidance mechanism, to make it more robust.

### 2.3 Learning with Label Noise and Face Uncertainties

Noisy labels in training dataset may hamper the performance of trained CNNs and thus lots of works focus on learning with label noise. Some researches proposed robust loss functions [26], [27], [28], [29] so that even noisy labels are fed for training models, they will not decrease the performance. Other works aim to model the noise in order to relabel, re-weight or remove them. Some methods [30], [31], [32], [33], [34] identify and correct suspicious labels to their corresponding true class. Some methods [35], [36], [37], [38] try to assign importance weights to training samples. Usually clean samples will have larger weight value to reduce the influence of noisy labels. Some works [39], [40], [41], [42] aim to guide the network to choose clean instances for updating parameters. The co-teaching learning paradigm was proposed in [42] that two neural networks are trained simultaneously to teach each other with potential clean samples. In the field of face recognition, [43] aims to remove potential noise and construct a clean dataset while [44] proposed a noise-robust loss function to tackle with incorrect identities.

Recently, the uncertainties in the FER task has been illustrated. Due to ambiguous facial expressions, low-quality

facial images, inconsistent annotations and noisy labels, the network tends to fit these samples and thus causing over-fitting. Zeng et al. [2] firstly take the inconsistent annotation problem among different datasets into FER and try to find the latent truth with multiple models. Wang et.al. [4] suppresses the uncertainties by sample reweighting and relabeling scheme. The network learns the uncertainty of samples with a self-attention mechanism and uncertain samples will have small importance weights and then be relabeled for training. Fan et.al. [45] estimate the uncertainty by comparing the distance between sample and its class center in the embedding space. Outliers will have small weight and then weighted softmax is utilized to suppress the impact of uncertain instances.

Except utilizing large amount of unlabeled data to tackle the shortage of training data, our Progressive Teacher also takes the uncertainties of FER into account and makes the network robust to noisy labels. Studies on the memory effect of deep neural network show that deep networks will learn clean and easy instances in the early stage of training [46]. Intuitively, the clean and easy samples have smaller loss values while noisy samples have larger ones. As training continues, they will gradually memorize and overfit those noisy labels. We adopt this training scheme in [42] to train a noise-robust model as well as using auxiliary unlabeled images.

### 3 PROGRESSIVE TEACHER

Our proposed Progressive Teacher is a typical teacher-student model, which means that the student model improves itself gradually with the guidance of teacher model. More specifically, the teacher model is the average of student model’s weights during training process and has better performance, the student model is promoted by optimizing the classification loss of labeled data and, in the meanwhile, learning to produce consistent outputs with the teacher model. As a semi-supervised algorithm, unlabeled data as well as labeled ones are used to compute the later consistency loss by penalizing the difference in the outputs of two models with the same image. Under the guidance of teacher model, the student model improves increasingly and finally achieves similar learning ability. This semi-supervised training strategy is proved to work well tackling the shortage of data. However, when there are label noise in dataset, which is inevitable in FER task, the student model will gradually fit those noisy samples when optimizing the classification loss and then the performance will degrade. Additionally, this passive effect will accumulate in the teacher model, resulting in the decline of its guidance ability. Consequently, the overall performance will drop. Therefore, our Progressive Teacher utilizes a unified framework which is shown in Fig 3, to tackle both the data shortage and noisy labels simultaneously. It abandons potential noisy samples automatically and progressively to make the student model learn pure useful knowledge. Compared to traditional teacher-student model like Mean Teacher [11], the teacher not only provide better logits and guides student to produce similar outputs, but also selects potential clean samples for student models to learn. Additionally, we use two pairs of teacher-student

models and adopt a cross-guidance mechanism. The two groups complement each other and boost recognition.

The two groups of teacher-student models share the same architecture but have different initialization. In each group, the student model (SNet) minimizes loss function and is optimized by stochastic gradient descent (SGD), the teacher model (TNet) derives from its student model by exponential average moving. Due to different initialization between two groups, they will feed each other with different samples in the early stage of training. Trained with different samples, the two student models differ and this variance will accumulate in their teachers. In this circle, we argue that the two groups complement each other and will learn better. Similar to Mean teacher, we regard the average model as teacher model which performs better and guides the student model. Differently, the teacher not only provides guidance, but also select clean samples for student to learn. When inputting samples into T Nets, those with smaller loss (cross-entropy loss) values are believed to be clean and then conveyed to the S Nets in the other group for further training. A cross-guidance mechanism is adopted to boost performance. For example, selected labeled clean samples from TNet-1 will be conveyed to SNet-2, then SNet-2 will compute the supervised classification loss of them and the unsupervised consistency loss with TNet1 to update network weights. Specifically, we use cross-entropy loss as classification loss and mean squared error (MSE) as consistency loss.

Given labeled data set  $L = \{(x_i, y_i)\}$  and unlabeled data set  $U = \{x_j\}$ , the teacher and student model use different data augmentations when inputting images. We denote the student model as  $f(x, \theta, \eta)$  and teacher model as  $f(x, \theta', \eta')$ .

For convenience, the softmax of their outputs are also denoted as the same formulation. In each iteration with mini-batch  $B$ , clean samples selected by TNet-1 and TNet-2 are denoted as  $B_1$  and  $B_2$  respectively. The ratio of selected clean samples is represented as  $R(t)$  and changes its value along iteration. For SNet-1, the supervised cross-entropy loss  $L_{s1}$  is denoted as following formulation:

$$L_{s1} = \frac{1}{|L \cap B_2|} \sum_{i \in L \cap B_2} \sum_{j=1}^C y_{i,j} \log f_j(x_i, \theta_1, \eta) \quad (1)$$

in which  $y_{i,j}$  is the  $j$ -th value of the ground truth label of the  $i$ -th sample and  $f_j(\cdot)$  represents the  $j$ -th output value of softmax of student model.  $C$  is the number of classes.

The unsupervised consistency loss is

$$L_{u1} = \frac{1}{|B|} \sum_{i \in B} \left\| f(x_i, \theta'_2, \eta') - f(x_i, \theta_1, \eta) \right\|^2 \quad (2)$$

Here *unsupervised* means that only images are utilized, without their labels. So we use all samples for computing unsupervised loss, which means that labeled data are also included. The overall loss for SNet-1 is

$$L_1 = L_{s1} + \omega(t)L_{u1} \quad (3)$$

in which  $\omega(t)$  is unsupervised weight ramp-up function. Because the teacher model has poor guidance in the beginning of training, we should give the unsupervised loss a

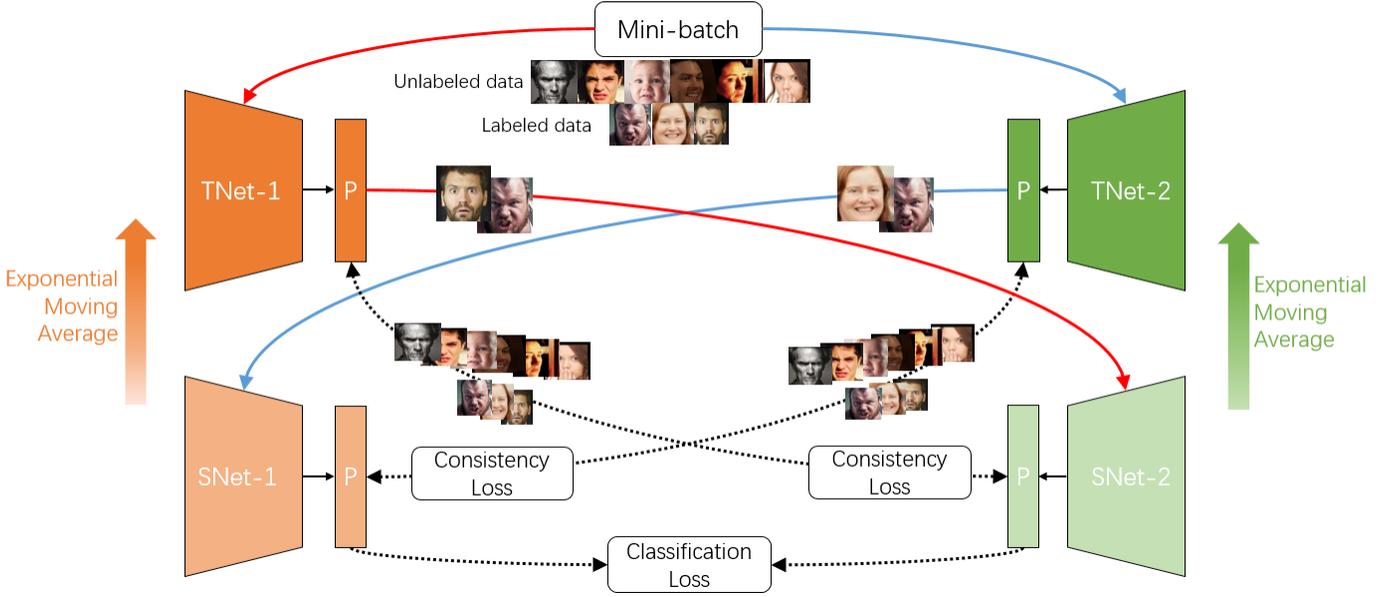


Fig. 3. The pipeline of our semi-supervised algorithm Progressive Teacher. Both labeled and unlabeled face images are fed into two teacher models with better learning capability (T Nets, TNet-1 and TNet-2) simultaneously. Then T Nets will select a part of labeled samples, which are believed to be useful and clean, and then fed them into the student models (S Nets, SNet-2 and SNet-1) in the other group. T Nets share the same network architecture with S Nets and derive from S Nets by exponentially averaging their weights along training steps. The number of selected clean samples decreases progressively to meet the training process and avoid remembering potential noisy samples. In the meanwhile, large-scale unlabeled data works that when feeding them into networks, S Nets improve themselves by learning to output consistent probability distributions with T Nets. For each SNet, it will compute the classification loss of selected labeled samples and the consistency loss using all samples in a mini-batch.

small weight first and increase it gradually. Similarly, the loss function for SNet-2 is denoted as

$$L_2 = L_{s2} + \omega(t)L_{u2} \quad (4)$$

in which

$$L_{s2} = \frac{1}{|L \cap B_1|} \sum_{i \in L \cap B_1} \sum_{j=1}^C y_{i,j} \log f_j(x_i, \theta_2, \eta) \quad (5)$$

and

$$L_{u2} = \frac{1}{|B|} \sum_{i \in B} \|f(x_i, \theta'_1, \eta') - f(x_i, \theta_2, \eta)\|^2 \quad (6)$$

In each iteration, after the student models update their weights, the teacher models are updated by following exponential moving average way.  $\alpha$  is a smoothing coefficient hyperparameter.

$$\theta'_t = \alpha \theta'_{t-1} + (1 - \alpha) \theta_t \quad (7)$$

As illustrated before, the network prefers to learn from clean samples and is not sensitive to noisy ones in the early stage of training, so we feed the S Nets with more labeled samples in the early stage and then progressively reduce the amount to avoid picking noisy ones. Algorithm 1 summarizes the learning process of our Progressive Teacher. We adopt two pairs of teacher-student models and use TNet from the other group to guide the SNet according to following considerations. Firstly, we use complementary T Nets to select clean samples instead of S Nets. Because S Nets are optimized with SGD, the parameters update quickly and are more affected by noisy labels. T Nets are more stable and learn better due to the accumulation along time.

---

#### Algorithm 1 The pipeline of Progressive Teacher.

---

##### Require:

The labeled data sets  $L$  and unlabeled data sets  $U$ ;  
 T Nets,  $f(x, \theta'_i, \eta')$ ,  $i = 1, 2$ ;  
 S Nets,  $f(x, \theta_i, \eta)$ ,  $i = 1, 2$

##### Ensure: Parameters $\theta_1, \theta_2, \theta'_1, \theta'_2$

Learning rate  $\lambda$ , cross-entropy loss  $l$ ,  $t = 0$

- 1: **for** epoch  $i \in [1, N]$  **do**
  - 2:   **for** mini-batch  $B$  **do**
  - 3:      $t = t + 1$
  - 4:     Select  $100 * R(t)\%$  small-loss samples by TNet-1:  
 $B_1 = \arg \min_{B': |B'| \geq R(t)|B|} l(f(x, \theta'_1, \eta'), B)$
  - 5:     Select  $100 * R(t)\%$  small-loss samples by TNet-2:  
 $B_2 = \arg \min_{B': |B'| \geq R(t)|B|} l(f(x, \theta'_2, \eta'), B)$
  - 6:     Update parameters in SNet-1:  
 $\theta_{1,t} = \theta_{1,t-1} - \lambda \frac{\partial L_1}{\partial \theta_{1,t-1}}$
  - 7:     Update parameters in SNet-2:  
 $\theta_{2,t} = \theta_{2,t-1} - \lambda \frac{\partial L_2}{\partial \theta_{2,t-1}}$
  - 8:     Update parameters in TNet-1:  
 $\theta'_{1,t} = \alpha \theta'_{1,t-1} + (1 - \alpha) \theta_{1,t}$
  - 9:     Update parameters in TNet-2:  
 $\theta'_{2,t} = \alpha \theta'_{2,t-1} + (1 - \alpha) \theta_{2,t}$
  - 10:    **end for**
  - 11: **end for**
- 

Therefore, small-loss samples selected by the later are more reliable. Secondly, we penalize the difference in outputs of SNet-1 (SNet-2) and TNet-2 (TNet-1). In semi-supervised learning, perturbation-based methods expect networks to produce consistent outputs when adding small noise to both

inputs and models themselves. To this end, we can get larger perturbation in models when using TNet and SNet from different groups. Additionally, SNet can learn more information from the variance of two groups since they are initialized differently and trained with not exactly the same samples. As the student model finally achieves comparable performance with its teacher, we choose one of the two teacher models for testing after training is finished.

## 4 EXPERIMENTS

In this section, we evaluate our method on RAF-DB and FERPlus with using AffectNet as auxiliary unlabeled data. Firstly, we will show the efficiency of semi-supervised learning. To further demonstrate the effect of tackling noisy labels, we add different levels of label noises to the above two datasets since the original datasets are believed to be reliable approximately. Specifically, for symmetric noise, we randomly select certain ratio of training samples in each expression category and change their labels to others uniformly. Considering expression surprised shares high similarity with fearful since they often occur with opened mouth or widened eyes, which may confuse annotators in real scenario, we also add asymmetric noise that the two expressions are relabeled to each other. Additionally, we conduct experiment on real-world noisy dataset AffectNet.

### 4.1 Datasets

**RAF-DB** [5] is a large-scale in-the-wild facial expression database which contains about 30,000 great diverse facial images from thousands of individuals downloaded from the Internet. Images in RAF-DB were labeled by 315 human coders and the final annotations were determined through the crowdsourcing techniques. And each image was assured to be labeled by about 40 independent labelers. RAF-DB contains 12,271 training samples and 3,068 test samples annotated with seven basic emotional categories (neutral expression is also taken into account). Compound expressions labeled as 11 classes are also provided, which are not used in our experiment.

**FERPlus** [7] is an extension of FER2013 [47]. The large-scale and unconstrained database FER2013 was created and labeled automatically by the Google image search API. All images in FER2013 have been registered and resized to 48\*48 pixels. FER2013 contains 28,709 training images, 3,589 validation images and 3,589 test images with seven expression labels. It's relabeled in 2016 by Microsoft with each image labeled by 10 individuals to consist 8 classes (contempt is added), thus has more reliable annotations. We only use seven basic emotional categories as RAF-DB.

**AffectNet** [6] contains more than 0.4 millions of labeled images which also includes contempt. The images are downloaded from Internet using three search engines and expression-related keywords. It's the largest dataset for FER currently. We choose the labeled images in 7-class basic expression categories, resulting in about 280,000 samples. In later experiment, we use these large amount of images as auxiliary unlabeled data without using their label information.

### 4.2 Implementation

For all images in RAF-DB, FERPlus and AffectNet, we detect face regions and locate their landmark points using MTCNN, then align them by similarity transformation. All images are further resized to 112\*112 pixels, which are shown in Fig.4. We use ResNet-18 as our backbone network, which is pre-trained on face recognition dataset CASIA-WebFace [48]. The dimension of feature is fixed to 512. All experiments are conducted under Pytorch framework, using one NVIDIA TITAN Xp GPU.



Fig. 4. Aligned images in three datasets.

For Progressive Teacher, the student model will improve from fast to slow, so we should gradually increase the value of  $\alpha$ . In the early stage of training,  $\alpha$  with small value makes the teacher model quickly forget previous inaccurate student weights. As training goes further, the teacher model benefits from long-term memory with larger  $\alpha$ . Specifically,  $\alpha = \min\{1 - 1/(1 + iter), 0.999\}$ , in which iter stands for iterations. We train the model for  $N$  epochs totally and  $N$  equals to 6 in experiment. The ratio of selected clean samples by teacher models

$$R(t) = 1 - r * \min\{t/T, 1\} \quad (8)$$

, in which  $r$  is set to 0.05 on original RAF-DB and FERPlus,  $r' + 0.1$  respectively when adding symmetric noise ratio of  $r'$ . When adding asymmetric noise to RAF-DB,  $r$  is set to 0.05 when noise ratio is 10%,  $r$  equals to 10% on other ratios. When adding asymmetric noise to FERPlus,  $r$  is set to 0.1 when noise ratio is 10%, 20%, 30%,  $r$  equals to 0.15 on other ratios.  $t$  represents the current iteration and  $T$  is the turning iteration and set to 300.  $R(t)$  decrease progressively to fit the training progress of neural networks that they are able to filter out noises automatically in the beginning and gradually overfit. After  $T$  iterations,  $R(t)$  keeps the same. The minibatch size is fixed to 200, in which 1/4 are labeled data and others are unlabeled. This facilitates the converge process and makes it controllable. We use random horizontal flip as data augmentation and SGD with momentum as optimizer. The momentum is 0.9 and weight decay is 0.0005. The learning rate starts with 0.01 and decreases to 0.001 after 3 epochs. The unsupervised weight ramp-up function  $\omega(t) = 10 * \exp\{-5 * (1 - epoch/N)^2\}$ .

### 4.3 Results

#### 4.3.1 Evaluations on reliable FER databases.

To demonstrate the effectiveness of this semi-supervised approach, we compare our algorithm with five other meth-

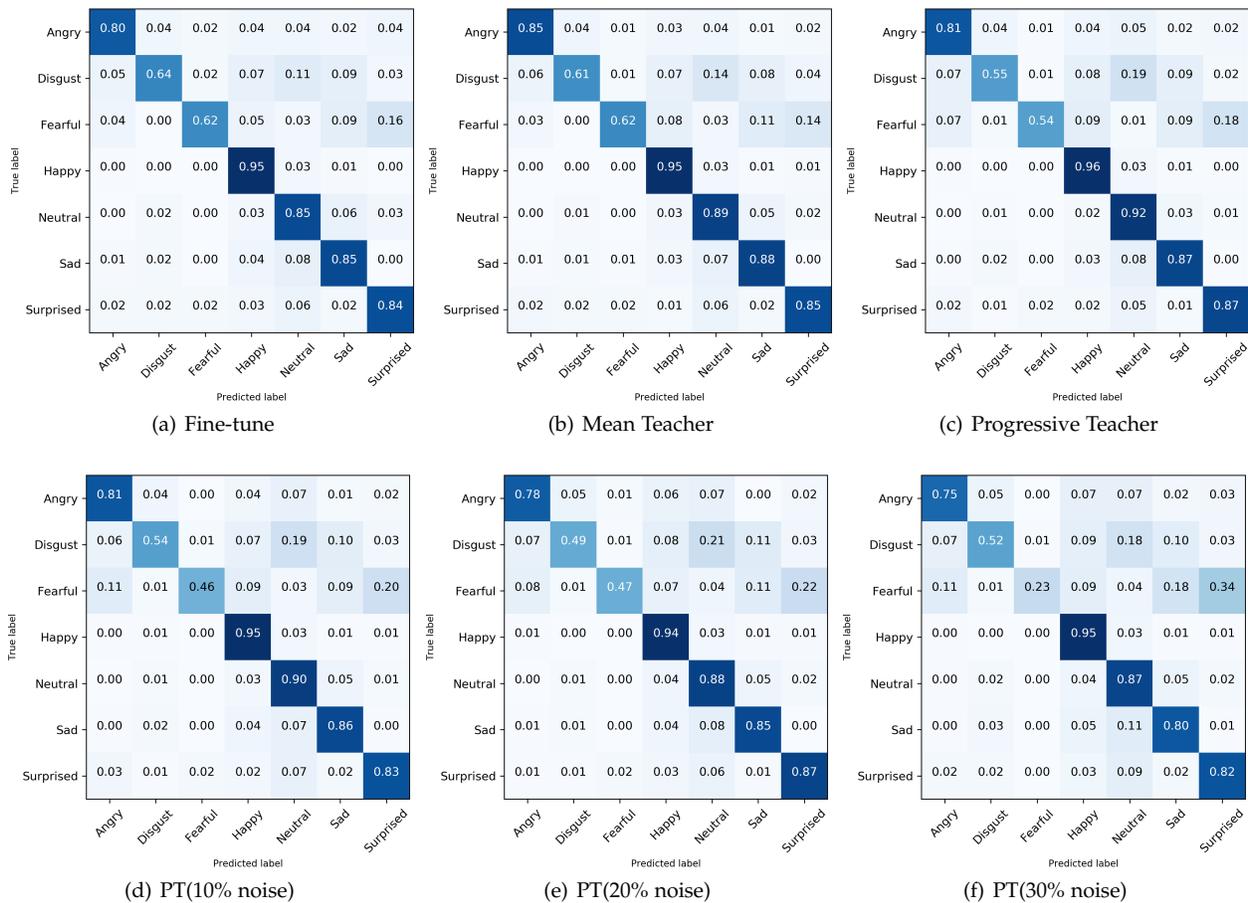


Fig. 5. The confusion matrices on RAF-DB. The top row represents results of three methods on original RAF-DB and the second row represents the results of Progressive Teacher under 3 levels of label noises on RAF-DB test set.

ods. The fine-tuning baseline means refining the pre-trained network on target dataset (RAF-DB or FERPlus). MixTrain uses all images and their labels in target dataset as well as the auxiliary dataset for training. In PseudoLabel method, we first train a network using target dataset, then classify samples in AffectNet to relabel them. These generated pseudo labels and corresponding images are thereby used to enlarge the target training set. It's worth mentioning that PseudoLabel is a simple semi-supervised method because it doesn't use the original labels of auxiliary data. State-of-art methods in semi-supervised learning, Mean Teacher [11] and FixMatch [49], are also provided for comparison. The results are shown in Table 2.

TABLE 2

Compare our method with baselines and semi-supervised learning method. ResNet-18 is used as backbone. †AffectNet is used as auxiliary dataset. ‡For RAF-DB, FERPlus is used as auxiliary dataset and vice versa.

Method	RAF-DB	FERPlus
Fine-tune(baseline)	87.09%	86.06%
MixTrain†	85.69%	85.40%
MixTrain‡	87.32%	85.85%
PseudoLabel	87.39%	85.13%
Mean Teacher [11]	88.41%	86.15%
FixMatch [49]	87.74%	86.45%
PT(ours)	<b>88.69%</b>	<b>86.60%</b>

Observing that the results of MixTrain are lower than baseline even with the huge amount of labeled images in AffectNet, this may be caused by two reasons. One is the annotation bias in different datasets, which is in accordance with [2], and the other is the existence of noisy labels in AffectNet. When utilizing two clean datasets which have rare noisy labels, the results of MixTrain also declare that enlarging data amount by merging multiple datasets can't improve remarkably and even degrade the performance due to inconsistent annotations. Compared to baseline, PseudoLabel improves the accuracy by 0.3% on RAF-DB but drops by 0.73% on FERPlus. The generated pseudo labels are not reliable enough due to inaccurate classifiers and domain shift. Images in RAF-DB shares more similarity with AffectNet so that the performance gets better due to the increase of training data. However, the large domain shift makes this semi-supervised method not work on FERPlus. Without utilizing label information of AffectNet, our Progressive Teacher improves the fine-tuning baselines by 1.6% on RAF-DB, and 0.54% on FERPlus respectively. Similar to the results of PseudoLabel, greater data similarity makes greater improvement on RAF-DB. Table 3 shows state-of-the-art results on RAF-DB. Compared with SCN, we both use pre-trained ResNet-18 as backbone and AffectNet as auxiliary data, but we don't employ its label information. This indicates the effectiveness of our semi-supervised ap-

TABLE 3  
Comparison to the state-of-the-art results on RAF-DB.

Method	Accuracy
PAT-ResNet-(gender, race) [50]	84.19%
IPA2LT(LTNet) [2]	86.77%
Acharya et al. [51]	87.0%
Gan et al. [52]	86.31%
APM-VGG [53]	85.17%
SCN [4]	88.14%
DAS + VGG-F [13]	86.55%
SCAN-CCI [54]	89.02%
PT-ResNet18 (ours)	88.69%
PT-SEResNet50-IR (ours)	<b>89.57%</b>

proach. Using SEResNet50-IR as backbone, we achieve the best accuracy of 89.57% on RAF-DB.

#### 4.3.2 Evaluations on synthetic noisy FER databases.

**Uniform Noise** To correspond with SCN [4], which also considers the situation of label noise, we randomly choose 10%, 20%, and 30% of training data for each category and randomly change their labels to others.

TABLE 4

The evaluation of PT on synthetic noisy RAF-DB and FERPlus. In SCN, "x" represent fine-tuning the pre-trained ResNet-18, "✓" means SCN algorithm is used. SCN doesn't use auxiliary datasets. †Pretrained SEResNet50-IR is used as backbone.

Method	Noise(%)	RAF-DB	FERPlus
SCN(x)	0	84.20%	86.80%
SCN(✓) [4]	0	87.03%	88.01%
RW Loss [45]	0	87.97%	87.60%
SCAN-CCI [54]	0	89.02%	82.35%
Mean Teacher [11]	0	88.41%	86.15%
Fine-tune(baseline)	0	87.09%	86.06%
PT(ours)	0	<b>88.69%/89.57%</b> †	86.60%
SCN(x)	10	80.81%	83.39%
SCN(✓) [4]	10	82.18%	84.28%
RW Loss [45]	10	82.43%	83.93%
SCAN-CCI [54]	10	84.09%	79.25%
Mean Teacher [11]	10	84.68%	82.87%
Fine-tune(baseline)	10	82.82%	82.09%
PT(ours)	10	<b>87.28%</b>	<b>85.07%</b>
SCN(x)	20	78.18%	82.24%
SCN(✓) [4]	20	80.80%	83.17%
RW Loss [45]	20	80.41%	83.55%
SCAN-CCI [54]	20	78.72%	72.93%
Mean Teacher [11]	20	81.40%	82.87%
Fine-tune(baseline)	20	77.97%	77.97%
PT(ours)	20	<b>86.25%</b>	<b>84.27%</b>
SCN(x)	30	75.26%	79.34%
SCN(✓) [4]	30	77.46%	82.47%
RW Loss [45]	30	76.77%	82.75%
SCAN-CCI [54]	30	70.99%	68.90%
Mean Teacher [11]	30	75.00%	72.06%
Fine-tune(baseline)	30	71.50%	70.03%
PT(ours)	30	<b>84.32%</b>	<b>83.73%</b>

In Table 4, we evaluate our PT algorithm under different levels of label noises on RAF-DB and FERPlus to demonstrate its effectiveness. We report the result when the network becomes stable instead of selecting the highest accuracy during training process. The results of state-of-the-art supervised methods on FER [4], [45], [54] under different levels of noise are also shown to make comparison. The backbones of SCN [4] and SCAN-CCI [54] are ResNet-18 and ResNet-50, respectively. We can see that on original

RAF-DB, Progressive Teacher surpasses Mean Teacher by only 0.28%, because the dataset contains rare noise. Under the noise rate of 10%, 20% and 30% on RAF-DB, our method all achieves the highest accuracy and improves the baseline by 4.46%, 8.28% and 12.82% respectively. The effect of our PT algorithm gets more remarkable as the noisy labels increase. Results on FERPlus show the same regularity. When the noise rate variates from 10% to 30%, the accuracy improvements are 2.98%, 6.30% and 13.7%. Though we get inferior accuracy on original FERPlus due to the domain discrepancy between FERPlus and AffectNet, we obtain superior performance on synthetic noisy datasets which indicates the advantage of semi-supervised approach and sample selection. Additionally, even adding 30% noises to original RAF-DB and FERPlus, the performance drops by only 4.37% and 2.87% respectively with the help of PT while SCN degrades by 9.57% and 5.54%. This indicates that our method can work in extremely hard condition. To see the performance's changing trend, we add more detailed noise levels on RAF-DB and the result is shown in Fig.6.

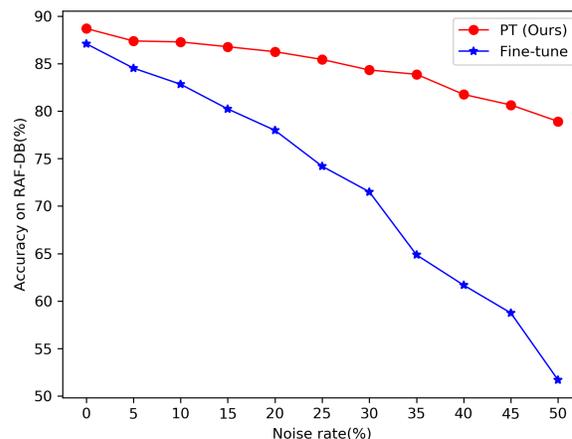
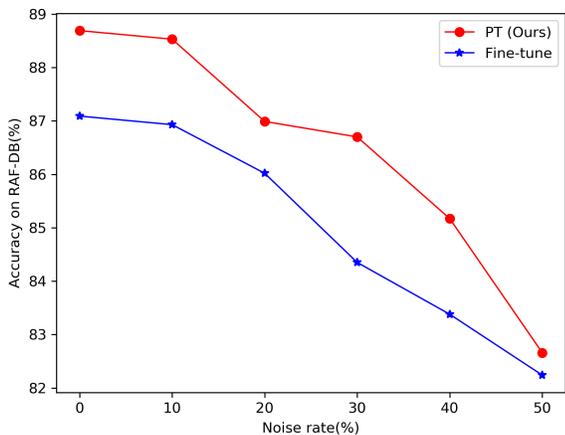


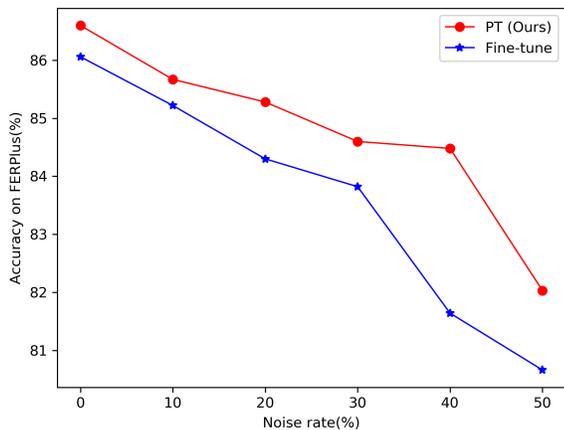
Fig. 6. The performance on RAF-DB under detailed noise levels.

Fig.5 shows the confusion matrices on RAF-DB. We can find that disgust and fearful are the hardest two expressions to recognize, which can be explained that they contain least samples. The class imbalance problem leads to the variance of recognition performance between different expressions. Our Progressive Teacher improves the overall accuracy but degrades on difficult expressions such as disgust and fearful. This is the natural consequence of data selection since samples in small-class are relatively hard to recognize and will have larger loss values so that some of them are likely to be filtered out by network. We focus on this problem and extend our method in section 4.6.

**Asymmetric Noise** In real situation, expression surprised is easy to be annotated as expression fearful and vice versa, we also conduct experiments under asymmetric noise. Specifically, samples in class fearful and class surprised are annotated to each other in the same ratio. The experimental result is recorded in Fig. 7. We can see that our method also achieves better performance than fine-tuning baseline.



(a) Performances on RAF-DB under asymmetric noise.



(b) Performances on FERPlus under asymmetric noise.

Fig. 7. The performance of Progressive Teacher when adding asymmetric noise to RAF-DB and FERPlus. In this scenario, fearful and surprised are relabeled to each other.

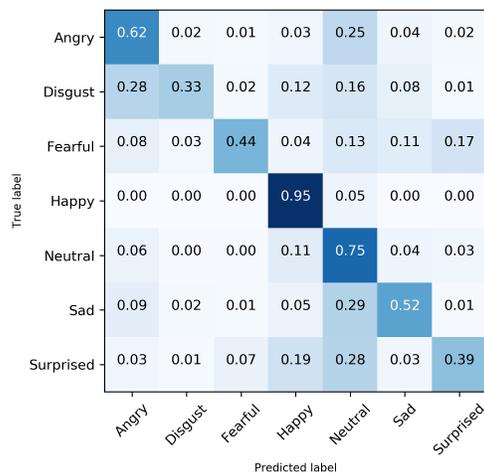
### 4.3.3 Evaluations on real-world noisy FER database Affect-Net.

We also conduct experiments on AffectNet, a real noisy database. FERPlus dataset is utilized as unlabeled data. The size of mini-batch is set to 200 and half of them are labeled data. The abandoning rate  $r$  is set to 0.2,  $T$  is set to 3000. Besides, among selected large-loss samples, disgust, fearful and surprised are reserved with the probability of 0.9, angry and sad are reserved with the probability of 0.8. We compare our method with fine-tuning baseline and Mean Teacher in Table. 5. The confusion matrices on AffectNet test set is

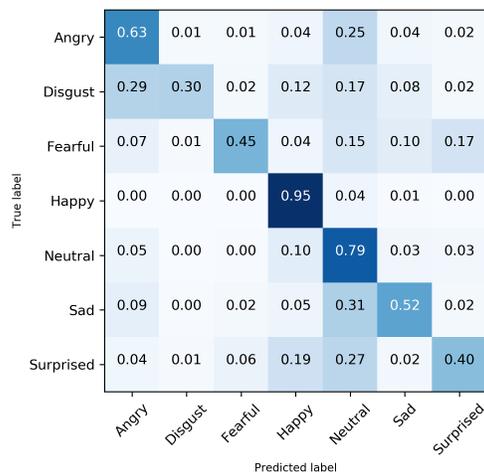
TABLE 5  
The accuracy when training with AffectNet.

	Test on AffectNet	Test on RAF-DB
Fine-tune	57.23%	79.04%
Mean Teacher [11]	57.62%	81.19%
PT(ours)	58.54%	82.11%

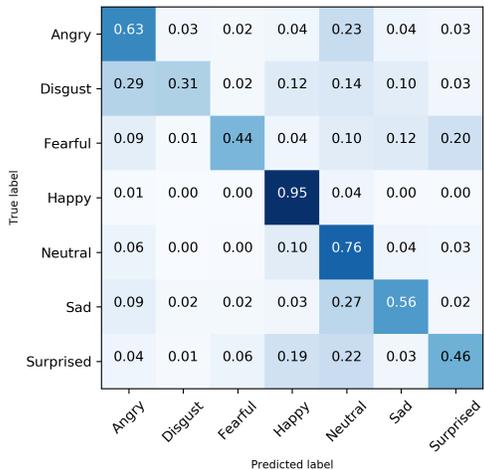
shown in Fig.8.



(a) Fine-tune



(b) Mean Teacher



(c) Progressive Teacher

Fig. 8. The confusion matrices on AffectNet.

## 4.4 Visualizations

To visualize the effect of PT algorithm, we show some examples which are filtered out as noisy samples in Fig.9.

This experiment is conducted on AffectNet. The selected training samples are with the largest loss values in each mini-batch. We also list the data distribution of samples in Table 6.

We add synthetic noise to RAF-DB since it contains rare noisy labels. Fig.10 and Table 7 shows the data distribution of abandoned samples under different noise levels on RAF-DB. Note that we abandon  $r'+10\%$  large-loss samples in mini-batch, in which the synthetic noise ratio is  $r'$ . The majority of abandoned samples are synthetic noisy labels and the rest are with original label.

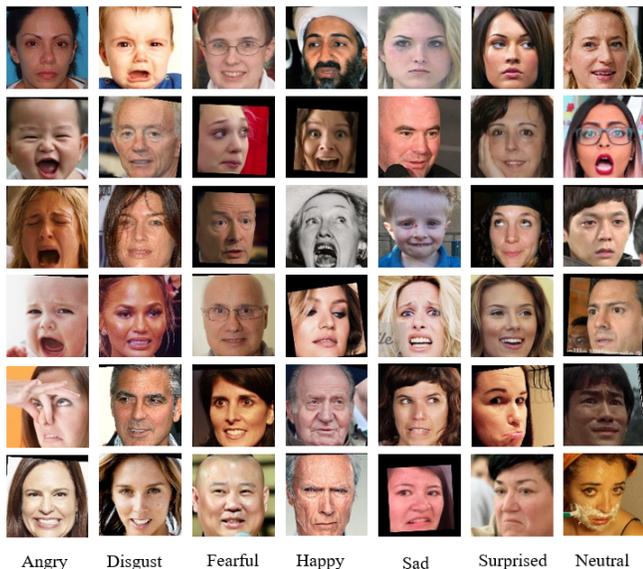


Fig. 9. Filtered training samples with Progressive Teacher.

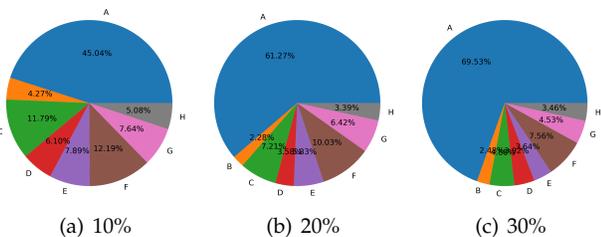


Fig. 10. The distribution of filtered samples in (A)synthetic noisy labels, (B)Angry, (C)Disgust, (D)Fearful, (E)Happy, (F)Neutral, (G)Sad, (F)Surprised when adding noise to RAF-DB.

TABLE 6  
The number of samples in each expression class of original datasets and filtered out data.

	Ang	Dis	Fea	Hap	Neu	Sad	Sur
RAF-DB	705	717	281	4772	2524	1982	1290
Filtered	49	149	96	52	46	63	37
FERPlus	2399	175	648	7410	9365	3403	3378
Filtered	102	72	79	86	347	297	89
AffectNet	24882	3803	6378	134415	74874	25459	14090
Filtered	1548	219	346	9078	19180	1843	592

TABLE 7  
The distribution of filtered out images on synthetic noisy RAF-DB.

Noisy rate	Synthetic Noise	Ang	Dis	Fea	Hap	Neu	Sad	Sur
10%	1108	105	290	150	194	300	188	125
20%	2261	84	266	132	215	370	237	125
30%	3421	122	240	193	179	372	223	170

### 4.5 Ablation Study

Progressive Teacher selects samples progressively which means that it utilizes less labeled data as training continues. In this subsection, we first conduct experiment to demonstrate the effectiveness of selecting samples progressively instead of keeping the same ratio all through, which is shown in Fig.11. The blue line represents using the same selecting ratio  $R = 1 - r$ . To correspond with the setting in Progressive Teacher, we set  $r$  to 0.2, 0.3 and 0.4 respectively as the noise rate increases. We can see that progressive selection outperforms using fixed selecting ratio under three noise rate levels. The noise rate gets higher, the former selection method gets more effective.

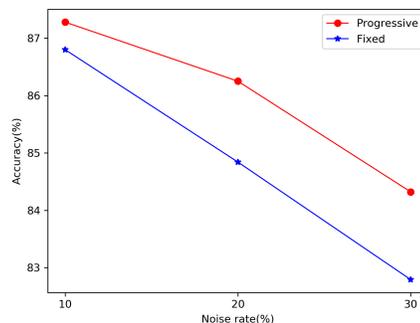


Fig. 11. The comparison between abandoning samples progressively and with the same ratio all through. Symmetric noise is adopted.

We also conduct experiments to illustrate the effectiveness of cross-guidance mechanism over same-guidance. When using one group of teacher-student model, which means that the teacher model feeds the student with clean samples and we compute the consistency loss between them, the accuracy of same-guidance mechanism on RAF-DB is 87.84%, which is lower than the 88.69% accuracy in cross-guidance setting by 0.85%. This indicates that the cross-guidance mechanism is superior to same-guidance. Additionally, considering there are two teacher-student groups but we compute the consistency loss between TNet-1 and SNet-1 (TNet-1 feeds SNet-2 with clean samples), we achieve accuracy of 88.10% on RAF-DB. When using two pairs of models, they have different weight initializations and thus will feed each other with different samples in the early stage of training. This will facilitate the variance of learning in two student models. And this variance will accumulate in their teachers. Different training samples make the two groups have different learning ability. In this circle, we argue that the two groups complement each other and will learn better. Moreover, adding more groups is difficult

to implement and brings more computing cost. Therefore, two pairs of networks is a good choice. We compare the performance of same-guidance and cross-guidance in Fig.12.

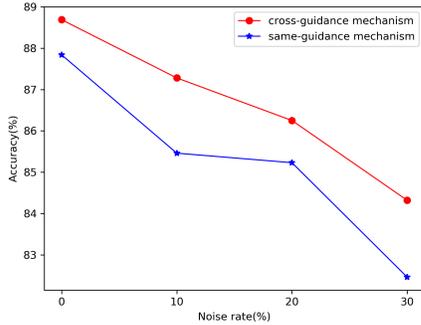


Fig. 12. The comparison between same-guidance and cross-guidance mechanism on RAF-DB. Symmetric noise is adopted.

#### 4.6 Further Study

In this section, we will analyze the relationship between loss value and noisy label first and then introduce an extension version of our method. Our experiments are conducted on RAF-DB. We count the average loss value of seven expressions and noisy labels and plot them in Fig.13. It can apparently be observed that noisy labels have much larger loss values than others. However, due to the lack of samples in some expressions (i.e., angry, disgust and fearful) and the difficulty of recognition, the loss value of these categories are larger than others (but smaller than noisy labels). In order not to improperly abandon these samples, we guide our Progressive Teacher with a confidence discriminator. Specifically, we train a confidence estimator which learns the confidence for the prediction of current sample. It's originally designed for out-of-distribution detection [55]. The trained confidence estimator has two heads, one of which outputs the prediction and the other outputs its confidence. Outliers and noisy labels usually have small confidence value. In our implementation, ResNet-18 is used as backbone and mean teacher mechanism is also adopted. When executing PT algorithm, if some large-loss samples in current mini-batch own large confidence for its prediction, these samples will be reserved and fed to student model. The reserved samples are intuitively clean but hard ones. The confidence threshold is set to 0.9. The accuracy of recognition and confusion matrices are shown in Fig.15. We can see that the recognition of expressions with much less samples (i.e., fearful) are better with the help of confidence estimator. We list some filtered out images on original RAF-DB in Fig.14. It shows good discernment on noisy labels and extreme uncertain samples.

### 5 CONCLUSIONS

In this paper, we propose a semi-supervised framework Progressive Teacher to tackle the shortage of data and inaccurate annotations in facial expression recognition so as to improve the recognition performance. The framework consists of two pairs of teacher-student models. In each

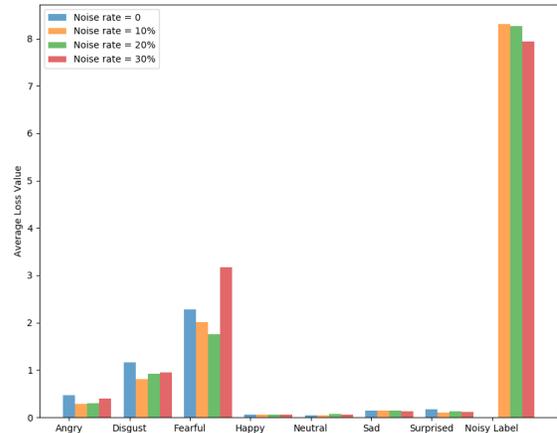
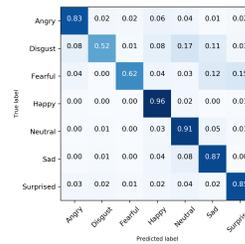


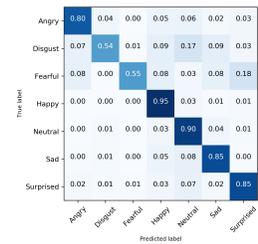
Fig. 13. The Average loss value of seven expressions and noisy labels on RAF-DB.



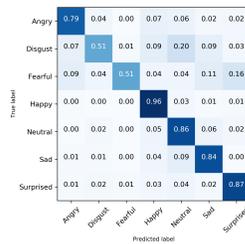
Fig. 14. Filtered training samples with Progressive Teacher and confidence estimator on original RAF-DB.



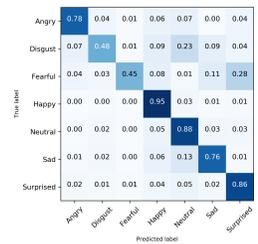
(a) No noise (Acc=88.59%)



(b) 10% noise (Acc=87.42%)



(c) 20% noise (Acc=86.47%)



(d) 30% noise (Acc=84.78%)

Fig. 15. The confusion matrices on RAF-DB when extending PT with confidence estimator.

pair, the student model computes supervised classification loss and unsupervised consistency loss and then update its parameters with SGD, the teacher model is the average of student model's weight during training process and guides its learning. Auxiliary large-scale unlabeled data is utilized to compute the unsupervised loss. Different from traditional semi-supervised learning method like Mean Teacher, our teachers can select potential clean samples for student models to learn and thus prevent from overfitting noisy samples. Additionally, we use the cross-guidance mechanism to boost performance. Extensive experiments show that our method achieves state-of-the-art result. We obtain the best recognition accuracy of 89.57% on RAF-DB.

## REFERENCES

- [1] 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)*. IEEE, 2017, pp. 118–126.
- [2] J. Zeng, S. Shan, and X. Chen, "Facial expression recognition with inconsistently annotated datasets," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 222–237.
- [3] S. Li and W. Deng, "Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 356–370, 2018.
- [4] K. Wang, X. Peng, J. Yang, S. Lu, and Y. Qiao, "Suppressing uncertainties for large-scale facial expression recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6897–6906.
- [5] S. Li, W. Deng, and J. Du, "Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2852–2861.
- [6] 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*, vol. 10, no. 1, pp. 18–31, 2017.
- [7] E. Barsoum, C. Zhang, C. C. Ferrer, and Z. Zhang, "Training deep networks for facial expression recognition with crowd-sourced label distribution," in *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, 2016, pp. 279–283.
- [8] 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.
- [9] M. Taini, G. Zhao, S. Z. Li, and M. Pietikainen, "Facial expression recognition from near-infrared video sequences," in *2008 19th International Conference on Pattern Recognition*. IEEE, 2008, pp. 1–4.
- [10] M. Valstar and M. Pantic, "Induced disgust, happiness and surprise: an addition to the mmi facial expression database," in *Proc. 3rd Intern. Workshop on EMOTION (satellite of LREC): Corpora for Research on Emotion and Affect*. Paris, France., 2010, p. 65.
- [11] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Advances in neural information processing systems*, 2017, pp. 1195–1204.
- [12] S. Li and W. Deng, "A deeper look at facial expression dataset bias," *IEEE Transactions on Affective Computing*, 2020.
- [13] P. Liu, Y. Wei, Z. Meng, W. Deng, J. T. Zhou, and Y. Yang, "Omni-supervised facial expression recognition: A simple baseline," *arXiv preprint arXiv:2005.08551*, 2020.
- [14] S. Li and W. Deng, "Deep facial expression recognition: A survey," *IEEE Transactions on Affective Computing*, 2020.
- [15] G. Levi and T. 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*, 2015, pp. 503–510.
- [16] X. Zhao, X. Liang, L. Liu, T. Li, Y. Han, N. Vasconcelos, and S. Yan, "Peak-piloted deep network for facial expression recognition," in *European conference on computer vision*. Springer, 2016, pp. 425–442.
- [17] S. Happy, A. Dantcheva, and F. Bremond, "A weakly supervised learning technique for classifying facial expressions," *Pattern Recognition Letters*, vol. 128, pp. 162–168, 2019.
- [18] M.-I. Georgescu and R. T. Ionescu, "Teacher-student training and triplet loss for facial expression recognition under occlusion," *arXiv preprint arXiv:2008.01003*, 2020.
- [19] J. E. Van Engelen and H. H. Hoos, "A survey on semi-supervised learning," *Machine Learning*, vol. 109, no. 2, pp. 373–440, 2020.
- [20] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Workshop on challenges in representation learning, ICML*, vol. 3, no. 2, 2013.
- [21] L.-Z. Guo, T. Han, and Y.-F. Li, "Robust semi-supervised representation learning for graph-structured data," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2019, pp. 131–143.
- [22] J. Du, C. X. Ling, and Z.-H. Zhou, "When does cotraining work in real data?" *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 5, pp. 788–799, 2010.
- [23] Z.-H. Zhou and M. Li, "Tri-training: Exploiting unlabeled data using three classifiers," *IEEE Transactions on knowledge and Data Engineering*, vol. 17, no. 11, pp. 1529–1541, 2005.
- [24] A. Rasmus, M. Berglund, M. Honkala, H. Valpola, and T. Raiko, "Semi-supervised learning with ladder networks," in *Advances in neural information processing systems*, 2015, pp. 3546–3554.
- [25] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," *arXiv preprint arXiv:1610.02242*, 2016.
- [26] A. Ghosh, N. Manwani, and P. Sastry, "Making risk minimization tolerant to label noise," *Neurocomputing*, vol. 160, pp. 93–107, 2015.
- [27] N. Charoenphakdee, J. Lee, and M. Sugiyama, "On symmetric losses for learning from corrupted labels," in *International Conference on Machine Learning*. PMLR, 2019, pp. 961–970.
- [28] A. Ghosh, H. Kumar, and P. Sastry, "Robust loss functions under label noise for deep neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.
- [29] Z. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," *arXiv preprint arXiv:1805.07836*, 2018.
- [30] A. Vahdat, "Toward robustness against label noise in training deep discriminative neural networks," *arXiv preprint arXiv:1706.00038*, 2017.
- [31] A. Veit, N. Alldrin, G. Chechik, I. Krasin, A. Gupta, and S. Belongie, "Learning from noisy large-scale datasets with minimal supervision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 839–847.
- [32] D. Tanaka, D. Ikami, T. Yamasaki, and K. Aizawa, "Joint optimization framework for learning with noisy labels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5552–5560.
- [33] K. Yi and J. Wu, "Probabilistic end-to-end noise correction for learning with noisy labels," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7017–7025.
- [34] E. Arazo, D. Ortego, P. Albert, N. O'Connor, and K. McGuinness, "Unsupervised label noise modeling and loss correction," in *International Conference on Machine Learning*. PMLR, 2019, pp. 312–321.
- [35] T. Liu and D. Tao, "Classification with noisy labels by importance reweighting," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 38, no. 3, pp. 447–461, 2015.
- [36] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *International Conference on Machine Learning*. PMLR, 2018, pp. 4334–4343.
- [37] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," *arXiv preprint arXiv:1902.07379*, 2019.
- [38] K.-H. Lee, X. He, L. Zhang, and L. Yang, "Cleannet: Transfer learning for scalable image classifier training with label noise," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5447–5456.
- [39] S. Reed, H. Lee, D. Anguelov, C. Szegedy, D. Erhan, and A. Rabinovich, "Training deep neural networks on noisy labels with bootstrapping," *arXiv preprint arXiv:1412.6596*, 2014.
- [40] H.-S. Chang, E. Learned-Miller, and A. McCallum, "Active bias: Training more accurate neural networks by emphasizing high variance samples," *arXiv preprint arXiv:1704.07433*, 2017.
- [41] E. Malach and S. Shalev-Shwartz, "Decoupling" when to update" from" how to update", *arXiv preprint arXiv:1706.02613*, 2017.

- [42] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *Advances in neural information processing systems*, 2018, pp. 8527–8537.
- [43] X. Wu, R. He, Z. Sun, and T. Tan, "A light cnn for deep face representation with noisy labels," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 11, pp. 2884–2896, 2018.
- [44] Y. Zhong, W. Deng, M. Wang, J. Hu, J. Peng, X. Tao, and Y. Huang, "Unequal-training for deep face recognition with long-tailed noisy data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7812–7821.
- [45] X. Fan, Z. Deng, K. Wang, X. Peng, and Y. Qiao, "Learning discriminative representation for facial expression recognition from uncertainties," in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 903–907.
- [46] D. Arpit, S. Jastrzebski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio *et al.*, "A closer look at memorization in deep networks," *arXiv preprint arXiv:1706.05394*, 2017.
- [47] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville, M. Mirza, B. Hammer, W. Cukierski, Y. Tang, D. Thaler, D.-H. Lee *et al.*, "Challenges in representation learning: A report on three machine learning contests," in *International conference on neural information processing*. Springer, 2013, pp. 117–124.
- [48] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning face representation from scratch," *arXiv preprint arXiv:1411.7923*, 2014.
- [49] K. Sohn, D. Berthelot, C.-L. Li, Z. Zhang, N. Carlini, E. D. Cubuk, A. Kurakin, H. Zhang, and C. Raffel, "Fixmatch: Simplifying semi-supervised learning with consistency and confidence," *arXiv preprint arXiv:2001.07685*, 2020.
- [50] J. Cai, Z. Meng, A. S. Khan, Z. Li, J. O'Reilly, and Y. Tong, "Probabilistic attribute tree in convolutional neural networks for facial expression recognition," *arXiv preprint arXiv:1812.07067*, 2018.
- [51] D. Acharya, Z. Huang, D. Pani Paudel, and L. Van Gool, "Covariance pooling for facial expression recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 367–374.
- [52] Y. Gan, J. Chen, and L. Xu, "Facial expression recognition boosted by soft label with a diverse ensemble," *Pattern Recognition Letters*, vol. 125, pp. 105–112, 2019.
- [53] Z. Li, S. Han, A. S. Khan, J. Cai, Z. Meng, J. O'Reilly, and Y. Tong, "Pooling map adaptation in convolutional neural network for facial expression recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 1108–1113.
- [54] D. Gera and S. Balasubramanian, "Landmark guidance independent spatio-channel attention and complementary context information based facial expression recognition," *Pattern Recognition Letters*, vol. 145, pp. 58–66, 2021.
- [55] T. DeVries and G. W. Taylor, "Learning confidence for out-of-distribution detection in neural networks," *arXiv preprint arXiv:1802.04865*, 2018.



**Weihong Deng** is a professor in School of Artificial Intelligence, Beijing University of Posts and Telecommunications. His research interests include computer vision and affective computing, with a particular emphasis in face recognition and expression analysis. He has published over 150 technical papers in international journals and conferences, such as IEEE TPAMI, TIP, IJCV, CVPR and ICCV. He serves as area chair for major international conferences such as IJ-CAI, ACM-MM, IJCB, FG, and ICME, and guest editor for IEEE TBIOM, and Image and Vision Computing Journal and the reviewer for dozens of international journals, such as IEEE TPAMI, TIP, TIFS, TNNLS, TAFCC, TMM, IJCV, and PR. His dissertation titled "Highly accurate face recognition algorithms" was awarded the Outstanding Doctoral Dissertation Award by Beijing Municipal Commission of Education in 2011. He has been supported by the program for New Century Excellent Talents in 2014, Beijing Nova in 2016, Young Chang Jiang Scholar, and Elsevier Highly Cited Chinese Researcher in 2020.



**Jing Jiang** received the B.E. degree in telecommunication engineering from the Beijing University of Posts and Telecommunications, Beijing, China, in 2020. She is currently pursuing the Ph.D. degree in information and telecommunications engineering. Her research interests include deep learning and facial expression analysis.