

A Training Details

A.1 The configuration of hyper-parameters

The training settings for ReCoT on three benchmarks are shown in Table 1. Most of the hyper-parameters are the same on the three benchmarks, except for total training epochs and epoch threshold K .

A.2 Details about baseline

Since existing noisy label learning (NLL) and semi-supervised learning (SSL) methods for multi-label AU recognition are rare, we re-implement some baselines with proper modification, e.g., mean-teacher (MT [5]), noise regularization (NR [10]), bootstrapping-hard (BS-H [4]), bootstrapping-soft (BS-S [4]), Mixup [6], knowledge distillation (KD [2]). For a fair comparison, all the backbones of these methods are ResNet34 with ImageNet pretrained initialization. For MT, we change its activation function to sigmoid to perform multi-label classification. The consistency loss of MT is still mean square error (MSE), which is the same as the original paper. For NR, we just change its classification loss, and the rest remains unchanged. Same as MT, we change the activation function of BS to the sigmoid. The linear weight of noisy labels for BS-H and BS-S is set to 0.8. The difference between BS-H and BS-S is the process of predictions. BS-S uses soft prediction probability to obtain fusion labels, while BS-H uses hard prediction labels. For Mixup, we perform the same operation as BS. For KD, we set the knowledge graph of different AUs as all 1's for simplicity. The linear weight of noisy labels for KD is set to 0.5.

B More Results

B.1 The Effectiveness of Noisy Data

Fig.1 presents the results of supervised, SSL, and WSL methods which are all implemented by us. Compared with the supervised methods, i.e., baseline, WSL methods such as NR and KD perform better in three databases, demonstrating that noisy data could provide useful information to improve performance. Apart from EmotioNet, the overall performance of WSL approaches is better than SSL methods, indicating there exists useful information in noisy labels.

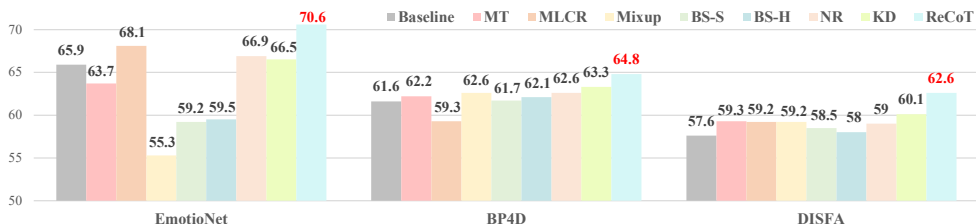


Figure 1: The comparison with baseline, MT [5], MLCR [4], Mixup [6], BS-S (H) [4], NR [10], KD [2] on three databases. Except for the results of MLCR and MT on EmotioNet, we re-implement all the other methods.

Table 1: The training configurations of hyper-parameters for ReCoT on three benchmarks. 046

Config	EmotioNet	BP4D	DISFA
optimizer	Adam	Adam	Adam
base learning rate	1e-3	1e-3	1e-3
batch size	100	100	100
total epochs	70	30	30
(α, β) for L_{noisy}	(1, 0.1)	(1, 0.1)	(1, 0.1)
selected ratio γ	0.9	0.9	0.9
fusion ratio μ	0.9	0.9	0.9
momentum ratio ϕ	0.9	0.9	0.9
epoch threshold K	40	15	15
balance coefficient λ_n	1	1	1
balance coefficient λ_{cons}	100	100	100

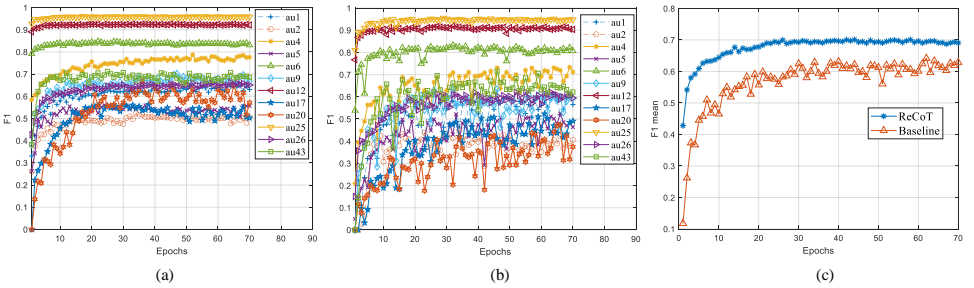


Figure 2: F1 scores by ReCoT and ResNet-34 are reported for different testing epochs on EmotioNet without ImageNet pre-training. (a) and (b) show the F1 scores for recognizing 12 AUs by ReCoT and ResNet-34, and (c) shows the average F1 scores by two methods. 070

B.2 Convergence Speed between Baseline and ReCoT 074

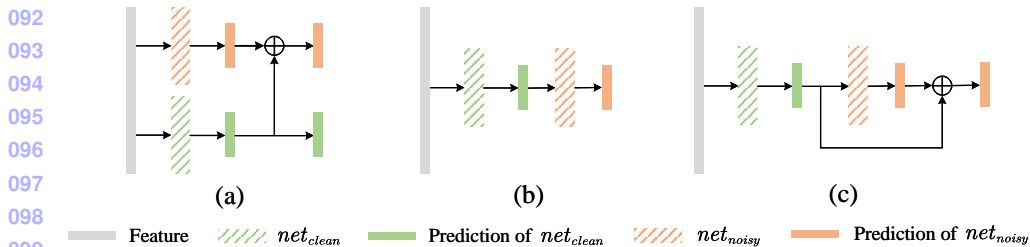
Fig. 2 shows that our ReCoT is useful for improving network convergence and the discriminative ability for AUs when exploiting the useful information from noisy data. Additionally, the training process of ReCoT is more stable compared with the baseline, since additional noisy data is utilized during the training process. 076

B.3 Different Noise modeling of net_{noisy} 081

We adopt three different types of label noise modeling for net_{noisy} (see Fig. (2)). Addictive label noise (Eq. 1) means modeling the label noise by net_{noisy} between noisy labels and clean predictions from net_{clean} . Multiplicative label noise (Eq. 2) indicates modeling the label noise transition matrix by net_{noisy} . Mixed label noise (Eq. 3) denotes that the label noise consists of additive and multiplicative label noise, which can be modeled using both the noise transition matrix and the residual connection. 083

$$p_{noisy} = net_{noisy}(x) + net_{clean}(x), \quad (1) \quad 090$$

$$p_{noisy} = net_{noisy}(net_{clean}(x)), \quad (2) \quad 091$$



100 Figure 3: Three types of label noise modeling of net_{noisy} . (a) Addictive label noise; (b)
101 multiplicative label noise; (c) mixed label noise.

$$102 \quad p_{noisy} = net_{noisy}(net_{clean}(x)) + net_{clean}(x). \quad (3)$$

103
104 Table 2 shows the results of three types of label noise modeling. From Table 2 we can
105 see that addictive noise modeling is good at modeling the label noise in the noisy set of
106 EmotioNet, and thus it achieves the best result compared with the other two methods.

107
108 Table 2: Three types of label noise modeling on EmotioNet.

Type of label noise modeling	Avg.
Addictive noise	70.6
Applicative noise	70.2
Mixed noise	70.3

109
110
111
112
113
114
115
116 Table 3: The choices of encoder under pretraining on EmotioNet. **Avg.** indicates the average
117 F1 scores of 12 AUs.

Type of Backbone	Avg.
ResNet18	65.2
ResNet34	65.9
ResNet50	65.3

118 119 120 121 122 123 124 125 126 B.4 The Influence of Different Encoder Choices

127 In the Ablation Study section of our main paper, we studied the influence when using differ-
128 ent encoders, such as ResNet-18, ResNet-34, and ResNet-50. The detailed results are given
129 in Table 3. We can notice that the F1 score of our method improves by 0.7% from ResNet-
130 18 to ResNet-34, but then drops by 0.6% from ResNet-34 to ResNet-50. A possible reason
131 is that the patterns related to AUs maybe not that complicated, and the overall modeling
132 difficulty comes from the label confusion for the AU recognition task.

133 134 135 B.5 The Influence of Different Noisy Image Scales

136 We briefly mentioned the influence of different noisy images Ablation Study section of our
137 main paper. Here, we show the detailed results in Table 4. It can be observed from Table 4

Table 4: The influence of the scale of noisy images on emotioNet. **Avg.** indicates the average F1 scores of 12 AUs.

Scale of Noisy Images	Avg.
10,000	69.8
50,000	70.6
100,000	70.2

that the performance improves 0.8% as the number of noisy images increases from 10,000 to 50,000, but gets saturated when further increasing the scale. We assume that 50,000 images can well represent the noisy label distribution in EmotioNet.

Table 5: The results of different views on EmotioNet. **Avg.** indicates the average F1 score of 12 AUs.

Views' Type	Avg.
View1	70.0
View2	70.1
The Average of Two Views	70.6

B.6 The Results of Different Views on EmotioNet

Since our ReCoT method contains two views, Table 5 shows the performance of two views and their average on EmotioNet during the inference phase. From Table 5, we can notice that the difference between two views is minor, while the average of two views is slightly better than a single view (about 0.5% higher). This suggests that when there is a low computation cost requirement, we may use a single view to perform inference, and the performance is still good; otherwise, we can use the average of two views for better performance.

184 References

- 185
- 186 [1] Mengying Hu, Hu Han, Shiguang Shan, and Xilin Chen. Weakly supervised image
187 classification through noise regularization. In *Proc. CVPR*, pages 11517–11525, 2019.
- 188 [2] Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li.
189 Learning from noisy labels with distillation. In *Proc. CVPR*, pages 1910–1918, 2017.
190
- 191 [3] Xuesong Niu, Hu Han, Shiguang Shan, and Xilin Chen. Multi-label co-regularization
192 for semi-supervised facial action unit recognition. In *Proc. NeurIPS*, pages 909–919,
193 2019.
- 194 [4] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and
195 Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping.
196 In *Proc. ICLR*, 2014.
197
- 198 [5] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-
199 averaged consistency targets improve semi-supervised deep learning results. In *Proc.*
200 *NeurIPS*, volume 30, 2017.
- 201 [6] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. Mixup: Be-
202 yond empirical risk minimization. In *Proc. ICLR*, 2018.
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229