

1 We thank all the reviewers for their valuable suggestions. Our response to individual reviewers' concerns are as follows.

2 =====To Reviewer 1=====

3 (1) **The differences between our work and [1]** include: (i) *The scope of the two papers is different.* While [1] is a
 4 fully supervised action unite (AU) recognition method, we focus on *semi-supervised* using massive unlabeled data
 5 and a small set of labeled data, which is more challenging but meaningful since labeling AU is difficult/expensive. (ii)
 6 The usage of AU relationship is different. While [1] used GNN to integrate the semantic relationship between AUs to
 7 enhance feature representation, in our work, leveraging GCN to encode AU relationship prior is only one part of our
 8 work, which can benefit the key part of our work, mining useful information from massive unlabeled data to obtain more
 9 informative and generalizable representation than learning from only a small labeled dataset. (iii) The performance on
 10 BP4D and DISFA. The performance of our semi-supervised learning method is lower than [1] on BP4D. However, we
 11 further conduct experiments on DISFA as [1] (with 100K unlabeled images from EmotioNet), and our method achieves
 12 56.8% Avg. F1 score, which is *higher* than that of [1] (55.9%). Although our semi-supervised learning method does
 13 not outperform the supervised learning method [1] on both datasets, we can still see the big potential of semi-supervised
 14 learning in AU recognition.

15 (2) **Consideration of mutually exclusive relations.** Our method also models the mutually exclusive relations of AUs.
 16 If two AUs are mutually exclusive, the avg. probability calculated by Eq. 7 will be small, and after normalization using
 17 Eq. 8, there will be a link added to the two AUs in the adjacent matrix.

18 (3) **Novelty of L_{mv} .** Our method learns two diverse classifiers in order to exploit diverse and informative features
 19 from unlabeled data for semi-supervised multi-label classification. Although the suggested CVPR19 paper also learns
 20 two diverse classifiers with paired labeled data for segmentation, the purpose is to predict how well each feature is
 21 semantically aligned between the source and target domains.

22 =====To Reviewer 2=====

23 (1) **Cross-database testing.** For cross-database testing, we train our model on EmotioNet with and without GCN
 24 and test the models on BP4D and UNBC (see Fig. 1(a)). Therefore, we agree that the AU co-occurrences may be
 25 different for different databases, but exploiting AU relationships provides better robustness of generated features for
 26 semi-supervised AU recognition under cross-database testing scenarios.

27 (2) **Language.** We will use copy-editing to improve our writing in the final version.

28 (3) **Justification for using two ResNet networks.** We conduct another experiment using ResNet-34 and Inception-v3
 29 Network, instead of using two ResNets. The avg. F1-score on EmotioNet is 67.6%, which is similar to using two
 30 ResNet-34 (68.1% F1 score). The results indicate that using two ResNets can generate features of two views which
 31 are different enough from each other. The main reasons are: (i) the two ResNets are initialized differently (pretrained
 32 separately); (ii) we have utilized L_{mv} to enforce them to generate different features.

33 (4) **Generalization to more views.** Currently, the L_{cr} and L_{mv} losses are designed for two views; one way to generalize
 34 to more views is to apply L_{cr} and L_{mv} to every two views. We will study this in future work.

35 (5) **Answers to the minor comments.** (a) PAC is a framework for mathematical analysis of machine learning, aiming
 36 at getting low generalization error with high probability. (b) "v" in Equ. 2 stands for "view". (c) We choose the number of
 37 unlabeled images according to the sizes of databases.

38 =====To Reviewer 3=====

39 (1) **Evidence of that two different networks can learn different cues for AU recognition.** We use t-SNE to visualize
 40 the features generated from the two views to recognize AU25 in Fig. 1(b). From the results, we can see that both views
 41 can achieve good classification accuracies, and the features generated from different views are very different, indicating
 42 that the two networks do learn different cues for AU recognition.

43 (2) **Explanation of the benefit of orthogonal weights.** Theoretically, the classifier with weights w and input feature f
 44 can be formulated as $\sigma(w^T f)$. After the model converges, the directions of vector w and f tend to be the same when
 45 the label is positive and tend to be opposite when the label is negative. So, w can be regarded as a representation of
 46 all the learned features. Therefore, orthogonalizing the weights will make the classifier weights independent to each
 47 other, and thus lead to the generated features from different views to be conditional independent because the feature
 48 generators and classifiers are optimized together. The feature visualization in Fig. 1(b) can also verify this conclusion.
 49 In addition, we also calculate the proportion of samples with inconsistent predictions from the two views with and
 50 without L_{mv} , and the results are shown in Fig. 1(c). From the results, we can see that orthogonalizing weights can
 51 make the predictions of the two views more different, and thus further benefit the semi-supervised co-training.

52 (3) **Fair comparison with the baseline model.** For fair comparisons, we use the *average F1 score* of the two ResNets
 53 without ensembling them as the final performance for both baseline and the proposed method, which guarantees that the
 54 proposed method is compared with the baseline under the same scale of parameters.

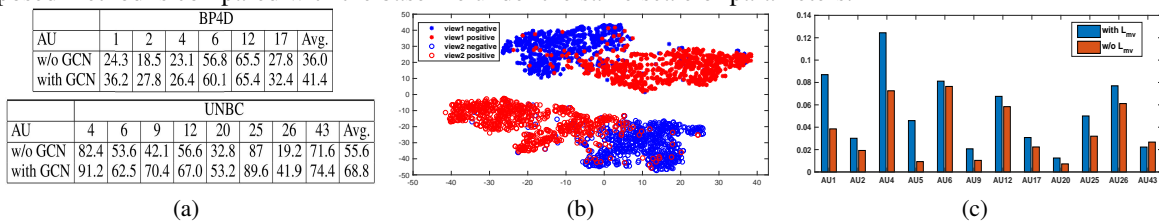


Figure 1: (a) F1 score (in %) for *cross-database testing* on BP4D and UNBC using models trained on EmotioNet. (b) t-SNE visualisation of the features generated from two views to recognize AU25. (c) Proportion of samples with inconsistent prediction results from the two views.

55

56 **Reference:** [1] Li G. et al. Semantic Relationships Guided Representation Learning for Facial Action Unit Recognition. AAAI2019.