

Supplementary Material for Multimodal Emotion Regression with Multi-Objective Optimization and VAD-Aware Audio Modeling for the 10th ABAW EMI Track

Anonymous CVPR submission

Paper ID 20

001 A. Additional Ablation Studies

002 A.1. Ablation on Temporal Enhancement Components

003
004 We further investigate whether more sophisticated temporal
005 enhancement modules can improve performance on top of
006 our pretrained multimodal feature pipeline. Starting from
007 the concatenation-based fusion setting, we separately intro-
008 duce three temporal components: Temporal Convolutional
009 Network (TCN), temporal BiGRU, and attention pooling.
010 The results are reported in Table 1.

011 The results show that *Baseline + TCN* achieves
012 0.466014, *Baseline + temporal BiGRU* achieves 0.463559,
013 and *Baseline + Attention Pooling* achieves 0.461779. None
014 of these variants surpass our final model. This obser-
015 vation suggests that, under our pretrained-feature regime,
016 the bottleneck does not lie in insufficient temporal mod-
017 eling. Instead, additional temporal enhancement may in-
018 troduce redundant transformations and optimization inter-
019 ference, thereby degrading the effectiveness of the original
020 pretrained representations.

021 A.2. Attempts on Textual Triplet Loss and Middle- 022 layer Prediction

023 In addition, we examine two alternative design choices be-
024 yond the main framework. The first is *Ours + tri loss*, which
025 follows prior work and introduces a text-guided triplet loss
026 for the textual modality. As shown in Table 2, this variant
027 yields 0.444169, substantially lower than the performance
028 of *Ours* (0.478567). This result indicates that the triplet-
029 based contrastive objective does not provide useful supervi-
030 sion for our continuous EMI regression task, and may even
031 distort the optimization target.

032 The second variant is *Ours + mid repr*, motivated by re-
033 cent studies suggesting that intermediate-layer representa-
034 tions may contain richer semantic information. In this set-
035 ting, intermediate features are directly used for prediction.
036 This variant achieves 0.472401, which is competitive but

Table 1. Ablation on temporal enhancement components on the validation set.

Method	Validation Score
Baseline + Temporal Enhancement	0.441066
Baseline + TCN	0.466014
Baseline + Temporal BiGRU	0.463559
Baseline + Attention Pooling	0.461779

Table 2. Attempts on textual triplet loss and middle-layer prediction on the validation set.

Method	Validation Score
Ours	0.478567
Ours + tri loss	0.444169
Ours + mid repr	0.472401

037 still below our final model. Therefore, while intermediate-
038 layer prediction preserves part of the discriminative infor-
039 mation, it does not lead to better performance than our de-
040 fault design.

041 Overall, these supplementary experiments further con-
042 firm that our final framework provides the best trade-off
043 between effectiveness and stability. The strongest perfor-
044 mance is obtained without extra temporal enhancement,
045 triplet supervision, or intermediate-layer prediction, sup-
046 porting our conclusion that, given strong pretrained fea-
047 tures, carefully designed fusion and optimization are more
048 important than adding extra temporal or auxiliary modules.

049 B. Comparison with Other Participating 050 Teams

051 We further compare our method with other participating
052 teams, as shown in Table 3. Our method obtains a test score
053 of 0.674, which is very close to *USTC-IAT-United* (0.680),
054 demonstrating competitive performance on the final evalu-

Table 3. Comparison with the results of the 8th participating teams.

Team	Validation Set	Test Set
HCAI-VIS	0.716	0.719
USTC-IAT-United	0.513	0.680
Ours	0.479	0.674

055 ation set.

056 While our validation score (0.479) is lower than those of
057 *HCAI-VIS* (0.716) and *USTC-IAT-United* (0.513), the much
058 smaller gap on the test set indicates that our model main-
059 tains good transferability to unseen data. In particular, the
060 difference between our method and *USTC-IAT-United* is
061 only 0.006 on the test set, showing that our approach re-
062 mains competitive under the final benchmark setting. Nev-
063 ertheless, the performance gap to *HCAI-VIS* suggests that
064 further improvements are still needed in terms of robustness
065 and ranking consistency across different data splits.