# Diagnosing and Re-learning for Balanced Multimodal Learning

Yake Wei[1], Siwei Li[2], Ruoxuan Feng[1], and Di Hu[✉1,3]

[1] Gaoling School of Artificial Intelligence, Renmin University of China, China
{yakewei,fengruoxuan,dihu}@ruc.edu.cn
[2] Department of Electronic Engineering, Tsinghua University, China
lisw19@mails.tsinghua.edu.cn
[3] Engineering Research Center of Next-Generation Search and Recommendation

**Abstract.** To overcome the imbalanced multimodal learning problem, where models prefer the training of specific modalities, existing methods propose to control the training of uni-modal encoders from different perspectives, taking the inter-modal performance discrepancy as the basis. However, the intrinsic limitation of modality capacity is ignored. The scarcely informative modalities can be recognized as "worse-learnt" ones, which could force the model to memorize more noise, counterproductively affecting the multimodal model ability. Moreover, the current modality modulation methods narrowly concentrate on selected worse-learnt modalities, even suppressing the training of others. Hence, it is essential to consider the intrinsic limitation of modality capacity and take all modalities into account during balancing. To this end, we propose the Diagnosing & Re-learning method. The learning state of each modality is firstly estimated based on the separability of its uni-modal representation space, and then used to softly re-initialize the corresponding uni-modal encoder. In this way, the over-emphasizing of scarcely informative modalities is avoided. In addition, encoders of worse-learnt modalities are enhanced, simultaneously avoiding the over-training of other modalities. Accordingly, multimodal learning is effectively balanced and enhanced. Experiments covering multiple types of modalities and multimodal frameworks demonstrate the superior performance of our simple-yet-effective method for balanced multimodal learning. The source code and dataset are available at https://github.com/GeWu-Lab/Diagnosing_Relearning_ECCV2024.

**Keywords:** Multimodal learning · Learning state diagnosing · Re-learning

## 1 Introduction

Inspired by the human's multi-sensory perception, multimodal learning where information from diverse sensors is jointly utilized, has witnessed tremendous progress in recent years [4, 11]. Even with current developments, the question

---
[✉]Corresponding author.

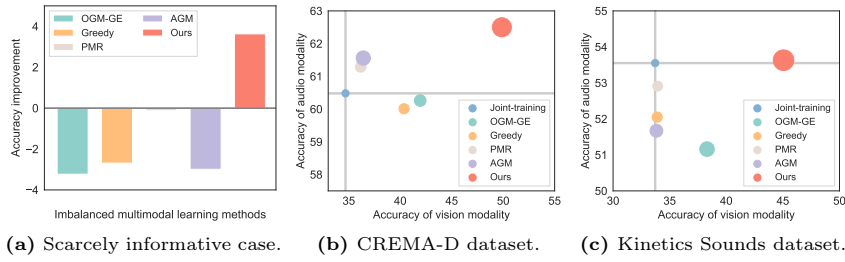**(a)** Scarcely informative case.      **(b)** CREMA-D dataset.      **(c)** Kinetics Sounds dataset.

**Fig. 1: (a): Scarcely informative modality case.** It shows the accuracy improvement compared with the joint-training baseline. Only our method has a positive performance improvement. **(b)&(c): Uni-modal encoder quality evaluation and comparison.** The uni-modal evaluation (Acc audio and Acc vision) is obtained by fine-tuning a new uni-modal classifier with the corresponding trained uni-modal encoder. A larger spot size reflects a better multimodal performance. Our method is superior in both multimodal performance and all uni-modal performance.

of how to assess and facilitate learning of individual modality remains open in the multimodal learning field. Especially, recent studies have found that some modalities in the multimodal model are less learnt than others [7, 15], called the imbalanced multimodal learning problem. This imbalance in modality utilization hinders the potential of multimodal learning, and even could make the multimodal model fail its uni-modal counterpart [15, 20]. Accordingly, a series of empirical methods are proposed to *balance* the uni-modal learning and achieve better multimodal performance [6, 9, 15, 25]. During the uni-modal balancing process, two keys are selection basis and balancing strategy.

In existing methods [6, 9, 15, 25], it is commonly believed that the modality with better prediction performance is the "well-learnt" modality, and correspondingly, the other "worse-learnt" modalities are the ones that need to be trained emphatically during uni-modal balancing. However, **they ignore the intrinsic limitation of modality capacity, where some modalities naturally have scarcely label-related information and more noise**. For cases of these modalities, the limited information causes their limited prediction performance, not just insufficient training. Although with worse prediction performance, purely emphasizing the training of these modalities could not bring many additional benefits and even force the model to memorize more noise, affecting the model ability. To further illustrate this problem, we modify the audio-vision CREMA-D dataset, and add white Gaussian noise into its audio data, making the audio modality with limited discriminative information but numerous noise. As Figure 1a, all existing imbalanced multimodal learning methods experience a performance drop compared with joint-training baseline. This phenomenon verifies that they wrongly push the training of scarcely informative modality with intrinsic limitation, counterproductively making them lose efficacy.

Upon the design of the balancing strategy, **existing methods narrowly concentrate on the learning of selected worse-learnt modalities [6, 9,**

**15,25].** Some even disturb the training of well-learnt modality [9,15], to facilitate the training of others. Inevitably, the ignorance or even suppression of well-learnt modality potentially affects its learning. As Figure 1b and Figure 1c, in existing imbalanced methods, although improving multimodal performance, the quality of the well-learnt audio modality can be worse than the joint-training baseline, especially on the Kinetics Sounds dataset. Recognizing these limitations, the challenge lies in how to overcome scarcely informative modalities cases and take all modalities into account during balancing.

In this paper, we propose the Diagnosing & Re-learning method, which periodically softly re-initialize the uni-modal encoder with representation separability as the basis. Since many multimodal models only have one multimodal output, it is hard to directly obtain the uni-modal learning state without additional modules. To this end, we focus on the uni-modal representation space, which is easier to access and can indirectly reflect the modality discriminative ability [17]. Concretely, the separability of train-representation and validation-representation are assessed by clustering, and utilized to diagnose the uni-modal learning state. In this way, the learning of each modality is well estimated individually. Then, to balance the uni-modal training, uni-modal encoders are softly re-initialized based on their learning state. **For well-learnt modalities**, their encoder has a greater re-initialization strength. This strategy helps the model reduce the reliance on them, and enhance the learning of other still under-fitting modalities. Simultaneously, the greater re-initialization of well-learnt modality avoids its over-training, even potentially improving the generalization [3, 31]. **For worse-learnt modalities**, their encoders are slightly re-initialized, which is also beneficial for escaping memorizing data noise that harms generalization. **When one modality is scarcely informative**, our method will not wrongly over-emphasize its training, and the re-initialization for its encoder helps encoders avoiding memorize data noise. Therefore, **our strategy can benefit all modalities at the same time**. In addition, the soft re-initialization partially preserves previously learnt knowledge already accrued by the network [3,18]. It safeguards the collaboration between modalities, ensuring that the collaborative knowledge is not completely discarded but rather fine-tuned.

Based on Figure 1a, our method can well handle the scarcely informative modality case and ideally achieves performance improvement. Moreover, as Figure 1b and Figure 1c, it also effectively enhances the learning of all modalities. Our method is flexible and can be equipped with diverse multimodal frameworks, including the multimodal Transformer. Overall, our contributions are three-fold. **Firstly,** we point out that existing imbalanced multimodal learning methods often ignore the intrinsic limitation of modality capacity and the well-learnt modality during balancing. **Secondly,** we propose the Diagnosing & Re-learning method to well balance uni-modal training by softly re-initialize encoders based on the uni-modal learning state. **Thirdly,** experiments across different types of modalities and multimodal frameworks substantiate the superior performance of our simple-yet-effective method.

## 2   Related Work

**Multimodal learning.** Motivated by the multi-sensory experiences of humans, the field of multimodal learning has gained significant attention and experienced rapid growth in recent years [11]. Multimodal learning involves the development of models capable of simultaneously integrating information from various modalities. Research in multimodal learning spans diverse domains, including such as multimodal recognition [26, 27] and audio-visual scene understanding [23, 32]. Commonly employed multimodal frameworks typically entail the extraction and fusion of uni-modal features, followed by the optimization of all modalities with joint learning objectives. However, besides the superficial multimodal performance across various tasks, the inherent learning of different modalities remains under-explored.

**Imbalanced multimodal learning.** Recent studies have revealed the imbalanced multimodal learning problem, where models prefer certain modalities over others, limiting their overall effectiveness [7, 15]. A variety of strategies have been suggested, focusing on balancing the optimization of individual modalities [6, 9, 21, 22, 25, 28]. For example, Peng *et al.* [15] proposed gradient modulation strategy, which dynamically monitors the contribution difference of various modalities to the final prediction during training and mitigates the gradient magnitude of dominant modality to focus more on other modalities. However, in these studies, they ignore the intrinsic limitation of modality capacity. One modality can be naturally scarcely informative and with plenty of noise. In addition, in these methods, when balancing uni-modal learning, only the selected worse-learnt modality is focused on. Some of them even intentionally impair the training of well-learned modalities to promote the learning of others [9, 15]. In this paper, we first diagnose the learning state of individual modality by its own representation separability without additional modules. Then, the learning state is used as the basis of soft encoder re-initialization to encourage further learning of under-fitting modalities while concurrently preventing the over-training of originally well-learnt modalities and scarcely informative modalities. It ensures better learning of all modalities.

**Network re-initialization.** Recent studies suggest that re-initializing the network parameters during training can effectively improve model performance and parameter utilization [1,3,18,31]. These methods involve re-initializing and transforming a part or all of the parameters of a neural network periodically. For instance, Alabdulmohsin *et al.* [1] proposed the Layer-wise Re-initialization strategy, which re-initializes the architecture block-by-block during training. Qiao *et al.* [16] proposed to detect unsatisfactory components in a neural network and re-initialize them to encourage them can better fit the tasks. Here we introduce the idea of network re-initialization, but with a very different intention. The uni-modal encoder is re-initialized based on its learning state, to ensure the benefits of all modalities as well as balancing uni-modal training.
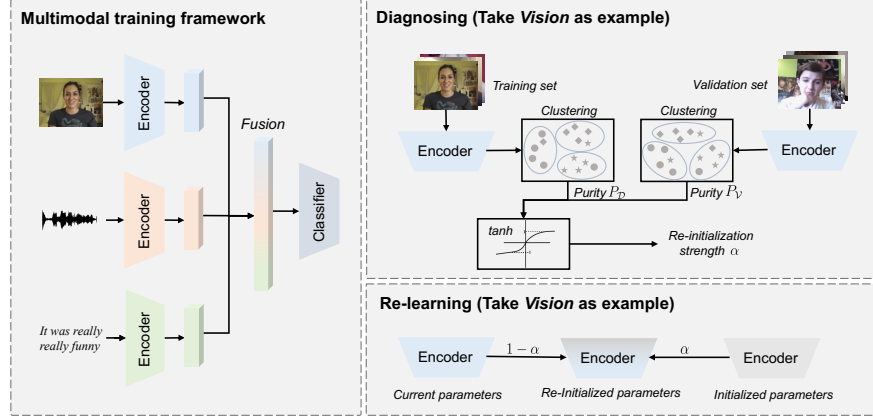
**Fig. 2:** Illustration of multimodal framework and the proposed Diagnosing & Re-learning method.

## 3   Method

### 3.1   Framework and notations.

**Multimodal framework.** As the left part of Figure 2, data of each modality is firstly fed into the corresponding uni-modal encoder to extract features. Then these uni-modal features are fused to obtain the multimodal feature. Our method has no reliance on the multimodal fusion strategy, and can cover simple fusion methods (*e.g.,* concatenation), and complex fusion methods (*e.g.,* cross-modal interaction). The fused feature is fed into the final multimodal classifier. One multimodal loss, cross-entropy, is utilized to optimize the model.

**Notations.** For the dataset with $K$ modalities, the training set is denoted as $\mathcal{D}$ with $N_{\mathcal{D}}$ samples and the validation set is denoted as $\mathcal{V}$ with $N_{\mathcal{V}}$ samples. Each data sample $x = \{x^1, x^2, \cdots, x^K\}$ is with $K$ modalities. The category number of the dataset is $M$. For each modality $k$, where $k \in \{1, 2, \cdots, K\}$, parameters of its encoder are denoted as $\theta_k$. $\theta_k^{\text{init}}$ represents the initialized parameter value.

### 3.2   Diagnosing: uni-modal learning state estimation

In multimodal learning, many multimodal models only have one multimodal output. Therefore, it is hard to directly obtain the uni-modal learning state without additional modules. In former studies, the estimation of uni-modal learning state often relies on specific fusion strategy [6, 15]. This limits their application to a wider range of scenarios. Elaborately designing ways to obtain uni-modal output is clearly complicated and not universal, since the multimodal fusion strategies are diverse. To well diagnose the uni-modal learning state without any additional modules or reliance on fusion strategies, we propose to focus on the uni-modal

representation space. It is known that the separability can reflect the representation quality [17]. Observing and comparing the separability of each extracted uni-modal representation is promising to capture the learning state. To evaluate representation separability, one straightforward idea is k-means clustering [13].

For data sample $x_i$ in the training set $\mathcal{D}$ with $N_{\mathcal{D}}$ samples, its $k-$th uni-modal feature that extracted by its $k-$th encoder $\theta_k$ is: $\phi_i^k = \theta_k(x_i^k)$. Then, to evaluate the separability of uni-modal features, it needs to split the set of all $k-$th uni-modal training features, $\Phi_{\mathcal{D}}^k = \{\phi_1^k, \phi_2^k, \cdots, \phi_{N_{\mathcal{D}}}^k\}$, into $M$ clusters. The set of all clusters is $\mathbf{C} = \{C_1, C_2, \cdots, C_M\}$, where $M$ is the category number.

Concretely, when splitting uni-modal features into clusters, $M$ samples in $\Phi_{\mathcal{D}}^k$ is firstly randomly picked as the centroid of $M$ clusters. Then, at the **assignment step**, each sample is assigned to the cluster with the nearest mean based on Euclidean distance. Concretely, sample $\phi_i^k$ is assigned to $m-$th cluster $C_m$ with centroid $O_m$ when:

$$\left\| \phi_i^k - O_m \right\|^2 \leq \left\| \phi_i^k - O_j \right\|^2 \quad \forall j, 1 \leq j \leq M. \tag{1}$$

$\| \cdot \|$ denotes the $L_2$-norm. After that, at the **updating step**, the centroid of each cluster is recalculated based on the current cluster:

$$O_m = \frac{1}{|C_m|} \sum_{\phi_i^k \in C_m} \phi_i^k. \tag{2}$$

After a given number of iterations between the assignment step and the updating step or the assignments no longer change, we have the final clustering results. For high-quality uni-modal representation, its ideal separability of feature space will bring satisfied clustering results. To evaluate the clustering results, we consider the clustering purity, which is a representative measurement for clustering quality [24]. Concretely, we first divide samples in $\Phi_{\mathcal{D}}^k$ into $M$ groups based on the ground truth labels and have classification sets $\mathbf{Z} = \{Z_1, Z_2, \cdots, Z_M\}$. Comparing the former clustering sets $\mathbf{C}$ and classification sets $\mathbf{Z}$, the purity is

$$P_{\mathcal{D}}^k = \frac{1}{N_{\mathcal{D}}} \sum_{C_m \in \mathbf{C}} \max_{Z_m \in \mathbf{Z}} |C_m \cap Z_m|. \tag{3}$$

It reflects the extent to which clusters contain a single class. Higher purity means better clustering results. And the uni-modal representation is of higher quality.

To diagnose the learning state of modality $k$, comparing the representation quality discrepancy between the training set $\mathcal{D}$ and the validation set $\mathcal{V}$ would be a useful reference. We know that when one model is well-learnt or even overtrained, its validation performance would be not increased according to the training performance [29]. This can also happen in the uni-modal encoder, bringing a gap between their train and validation representation quality. This gap is expected to reflect the learning state of one modality. Concretely, for the validation set $\mathcal{V}$, we also conduct the clustering algorithm and obtain its purity $P_{\mathcal{V}}^k$. And the gap between training set purity $P_{\mathcal{D}}^k$ and validation set purity $P_{\mathcal{V}}^k$ is:

$$g^k = |P_{\mathcal{D}}^k - P_{\mathcal{V}}^k|. \tag{4}$$

---

**Algorithm 1** Diagnosing & Re-learning

---

**Require:** Training set $\mathcal{D}$, validation set $\mathcal{V}$, epoch number $T$, uni-modal encoder parameters $\theta_k$, $k \in \{1, 2, \cdots, K\}$, Diagnosing & Re-learning frequency $H$.
  **for** $t = 0, \cdots, T - 1$ **do**
    Train and update parameters;
    **if** $t \mod H == 0$ **then**
      **for** $k = 1, \cdots, K$ **do**
        Extract training feature set $\Phi_{\mathcal{D}}^k$ and validation feature set $\Phi_{\mathcal{V}}^k$;
        Conduct clustering algorithm on $\Phi_{\mathcal{D}}^k$ and have its purity $P_{\mathcal{D}}^k$;
        Conduct clustering algorithm on $\Phi_{\mathcal{V}}^k$ and have its purity $P_{\mathcal{V}}^k$;
        Calculate the purity gap $g^k$ based on Equation 4;
        Calculate re-initialization strength $\alpha_k$ based on Equation 5;
        Reinitialize encoder parameters $\theta_k$ with $\alpha_k$ based on Equation 6.
      **end for**
    **end if**
  **end for**

---

Based on the property of purity, $g^k \in [0, 1]$. This purity gap reflects the quality gap between train and validation representation (Observations about this gap are provided in Section 4.8). When the value of purity gap $g^k$ is large, this modality is well-learnt or even over-trained. In this way, the learning state of one modality is diagnosed individually.

### 3.3 Re-learning: uni-modal re-initialization based on learning state

In Section 3.2, the uni-modal learning state is diagnosed by the separability discrepancy between training and validation representation space. Then, to balance the uni-modal training, we propose to softly re-initialize all uni-modal encoders based on their diagnosed learning state. This re-initialization breaks the model's reliance on one specific modality, and potentially enhances the model's generalization ability by re-learning multimodal data. Specifically, the re-initialization strength $\alpha_k$ for modality $k$ is calculated based on purity gap:

$$\alpha_k = \tanh(\lambda \cdot g^k), \tag{5}$$

where $\lambda > 1$ is the hyper-parameter to further control the re-initialization strength. Then we can have $\lambda \cdot g^k \geq 0$ and $\alpha_k \in [0, 1)$. The use of function $\tanh(x)$[4] aims to map the final re-initialization strength to a value between 0 and 1, while ensuring a monotonically increasing property when $x \geq 0$. These properties make the re-initialization strength $\alpha_k$ proportional to the purity gap $g^k$. Then, the encoder parameters of modality $k$ are re-initialized by:

$$\theta_k = (1 - \alpha_k) \cdot \theta_k^{\text{current}} + \alpha_k \cdot \theta_k^{\text{init}}, \tag{6}$$

---

[4] Other functions that satisfy these properties can also be used. We provide more experiments in the supplementary material.

where $\theta_k^{\text{current}}$ is the current parameter and $\theta_k^{\text{init}}$ is the initialized parameter.

With our strategy, on the one hand, for the well-learnt modalities, its encoder experiences a greater re-initialization, which effectively makes the model temporarily get rid of the dependence on them and enhances the learning of other still under-fitting modalities. Meanwhile, after re-initialization, the model would re-learn the former well-learnt data. This process can help to prevent the model from confidently fitting to the noise, avoiding over-training for well-learnt modalities. On the other hand, for other modalities (even they are scarcely informative), the slight re-initialization in their encoder also helps to prevent the memorization of data noise that negatively affects generalization. Overall, our method can benefit all modalities simultaneously. In addition, our soft re-initialization maintains a portion of the knowledge previously acquired by the model. It protects the learnt inter-modal correlation to some extent, making sure that the shared knowledge is not entirely lost but instead refined. The proposed method is illustrated in Figure 2, and the entire training process is shown in Algorithm 1. The Diagnosing & Re-learning strategy is conducted every $H$ epoch.

## 4 Experiment

### 4.1 Dataset

**CREMA-D** [5] is an emotion recognition dataset with two modalities, audio and vision. This dataset covers six emotions: angry, happy, sad, neutral, discarding, disgust and fear. The whole dataset contains 7442 clips.
**Kinetic Sounds** [2] is an action recognition dataset with two modalities, audio and vision. This dataset contains 31 human action classes, which are selected from the Kinetics dataset [8]. It contains 19k 10-second video clips.
**UCF-101** [19] is an action recognition dataset with two modalities, RGB and optical flow. This dataset contains 101 categories of human actions. The entire dataset is divided into a 9,537-sample training set and a 3,783-sample test set according to the original setting.
**CMU-MOSI** [30] is a sentiment analysis dataset with three modalities, audio, vision and text. It is annotated with utterance-level sentiment labels. This dataset consists of 93 movie review videos segmented into 2,199 utterances.

### 4.2 Experimental settings

For the CREMA-D and the Kinetic Sounds dataset, ResNet-18 is used as the backbone and models are trained from scratch. For the UCF-101 dataset, ResNet-18 is also used as the backbone and is ImageNet pre-trained. For the CMU-MOSI dataset, transformer-based networks are used as the backbone [10] and the model is trained from scratch. The choices of architecture and initialization follow former imbalanced multimodal learning studies, to have a fair comparison. During training, we use the SGD optimizer with momentum (0.9) and set the learning rate at $1e - 3$. All models are trained on 2 NVIDIA RTX 3090 (Ti).

**Table 1:** Comparison with imbalanced multimodal learning methods where bold and underline represent the best and second best respectively. Joint-training is the widely-used baseline with concatenation fusion and one multimodal loss function.

| Method | CREMA-D (Audio/Vision) | | Kinetics Sounds (Audio/Vision) | | UCF-101 (RGB/Optical Flow) | |
|---|---|---|---|---|---|---|
| | Acc | Macro F1 | Acc | Macro F1 | Acc | Macro F1 |
| Joint-training | 67.47 | 67.80 | 65.04 | 65.12 | 80.41 | 79.40 |
| G-Blending [20] | 69.89 | 70.41 | 68.60 | 68.64 | 81.73 | 80.84 |
| OGM-GE [15] | 68.95 | 69.39 | 67.15 | 66.93 | 81.15 | 80.36 |
| Greedy [25] | 68.37 | 68.46 | 65.72 | 65.80 | 80.60 | 79.50 |
| PMR [6] | 68.55 | 68.99 | 65.62 | 65.36 | 81.36 | 80.37 |
| AGM [9] | 70.16 | 70.67 | 66.50 | 66.49 | 81.55 | 80.58 |
| Ours | **75.13** | **76.00** | **69.10** | **69.39** | **82.11** | **80.87** |

**Table 2:** Additional comparison with G-Blending [20] method. For a fair comparison, we also introduce the same uni-modal classifier and uni-modal cross-entropy loss function as G-Blending in the Ours† method.

| Method | CREMA-D (Audio/Vision) | | Kinetics Sounds (Audio/Vision) | | UCF-101 (RGB/Optical Flow) | |
|---|---|---|---|---|---|---|
| | Acc | Macro F1 | Acc | Macro F1 | Acc | Macro F1 |
| Joint-training | 67.47 | 67.80 | 65.04 | 65.12 | 80.41 | 79.40 |
| G-Blending [20] | 69.89 | 70.41 | 68.60 | 68.64 | 81.73 | 80.84 |
| Ours | 75.13 | 76.00 | 69.10 | 69.39 | 82.11 | 80.87 |
| Ours† | **79.30** | **79.58** | **72.17** | **72.02** | **83.05** | **82.17** |

In experiments, our method is conducted every 5 epoch for the CMU-MOSI dataset and every 20 epoch for others. $\lambda$ is 3, 7, 7, 8 for CREMA-D, Kinetics Sounds, UCF-101 and CMU-MOSI dataset respectively. More ablation studies and comparisons are provided in the supplementary material.

### 4.3 Comparison with imbalanced multimodal learning methods

To assess the efficacy of our method in addressing the imbalanced multimodal learning problem, we conduct comparisons with recent studies: **G-Blending [20], OGM-GE [15], Greedy [25], PMR [6] and AGM [9].** Joint-training is the widely used baseline for the imbalanced multimodal learning problem, with concatenation fusion and one multimodal cross-entropy loss function [6, 9, 15]. The results are shown in Table 1. We provide experiments across several datasets with different modalities, like audio, vision and optical flow. Based on the results, we first observe that all these imbalanced multimodal learning methods achieve improvement in the multimodal performance, exhibiting the existence of the imbalanced multimodal learning problem and the necessity of balancing uni-modal learning during training. More than that, our method consistently exhibits superior performance across multiple datasets with different types of

**Table 3: (Left): Comparison with imbalanced multimodal learning methods on CMU-MOSI dataset with three modalities.** The greedy method could not extend to cases with more than two modalities. * indicates that the original methods of OGM-GE and PMR only consider two modality cases, but we extend them while retaining the core uni-modal balancing strategy. **(Right): Comparison with imbalanced multimodal learning methods with Transformer-based backbone on CREMA-D dataset.**

| Method | CMU-MOSI Transformer-based backbone (Audio/Vision/Text) | | | CREMA-D Transformer-based backbone (Audio/Vision) | | |
|---|---|---|---|---|---|---|
| | Acc | Weighted F1 | Macro F1 | Acc | Weighted F1 | Macro F1 |
| Joint-training | 76.96 | 76.76 | 75.68 | 68.55 | 68.88 | 69.33 |
| G-Blending [20] | 77.26 | 77.20 | 76.27 | 69.35 | 68.77 | 69.63 |
| OGM-GE [15] | 77.41* | 77.17* | 76.09* | 70.70 | 70.67 | 71.18 |
| Greedy [25] | / | / | / | 68.41 ($\downarrow$) | 68.72 ($\downarrow$) | 69.17 ($\downarrow$) |
| PMR [6] | 77.55* | 77.50* | 76.58* | 69.89 | 69.70 | 70.07 |
| AGM [9] | 77.26 | 77.07 | 76.02 | 69.22 | 69.34 | 69.61 |
| Ours | **77.99** | **78.09** | **77.37** | **71.64** | **71.66** | **72.12** |

modalities, outperforming other methods. This demonstrates the effectiveness of our Diagnosing and Re-learning strategy, which takes all modalities into account.

In Table 1, the G-Blending method also achieves considerable improvement in the model performance, especially on the Kinetics Sounds and UCF-101 datasets. However, compared with other methods that only have one multimodal joint loss, it introduces additional uni-modal classifier and correspondingly uni-modal loss functions. These additional modules are definitely helpful for controlling the training of individual modalities. As stated before, our method is flexible and not limited to special multimodal frameworks. Therefore, to have a fair comparison, we also introduce the same uni-modal classifier and loss function as G-Blending. The results are shown in Table 2. After introducing the same uni-modal classifier and loss function, our method (Ours† in the table) significantly outperforms the G-Blending method. In addition, these results also suggest that our method of targeted uni-modal encoder re-initialization on the basis of learning state can be effectively integrated with other modules while maintaining its effectiveness.

### 4.4 Comparison in more general multimodal frameworks

In multimodal learning, besides the widely used late-fusion framework with the convolutional neural network backbone, more complex transformer-based backbone and cross-modal interaction modules also have a wide range of applications. In this section, we conduct a comparison of the transformer-based multimodal frameworks. Results are shown in Table 3. For the CMU-MOSI dataset, Transformer is used as the backbone, following [10] and the model is trained from scratch. For the CREMA-D dataset, experiments in this section use the representative multimodal Transformer backbone, MBT [14]. It has both single-

**Table 4:** Comparison with imbalanced multimodal learning methods on scarcely informative modality case. We modify the audio data of the CREMA-D dataset, adding extra white Gaussian noise to make it noisier and scarcely discriminative. All compared methods have a clear performance drop in this case. In contrast, our method still achieves considerable improvement in this challenging scenario.

| Method | Scarcely informative modality case (Audio/Vision) | | |
|---|---|---|---|
| | Acc | Weighted F1 | Macro F1 |
| Joint-training | 65.86 | 66.03 | 66.39 |
| G-Blending [20] | 61.16 ($\downarrow$) | 60.45 ($\downarrow$) | 61.15 ($\downarrow$) |
| OGM-GE [15] | 62.63 ($\downarrow$) | 64.21 ($\downarrow$) | 65.07 ($\downarrow$) |
| Greedy [25] | 63.17 ($\downarrow$) | 62.99 ($\downarrow$) | 63.83 ($\downarrow$) |
| PMR [6] | 65.73 ($\downarrow$) | 64.81 ($\downarrow$) | 65.33 ($\downarrow$) |
| AGM [9] | 62.87 ($\downarrow$) | 62.28 ($\downarrow$) | 63.73 ($\downarrow$) |
| Ours | **69.49** | **69.72** | **70.28** |

modal layers and cross-modal interaction layers with dense cross-modal interaction modules. In experiments, the model is ImageNet pre-trained. The results lead us to make the following observation that existing imbalanced multimodal learning methods may lose efficacy in Transformer-based frameworks involving cross-modal interactions. For instance, the Greedy method [25] is even inferior to the joint-training baseline on the CREMA-D dataset. Conversely, our method, which has no reliance on special types of multimodal frameworks, demonstrates its ideal versatility and superior performance.

### 4.5    Comparison in more-than-two modality case

In current imbalanced multimodal learning methods, many of them only focus on the case of two modalities [6, 15, 25]. This limitation greatly hampers their application in broader scenarios. In contrast, our method has no restriction on the number of modalities. In this section, we compare these methods in the more-than-two modality case, on the CMU-MOSI dataset with three modalities: audio, vision and text. Among existing imbalanced multimodal methods, the Greedy [25] method could not be extended to this case suitably. In the original paper of OGM-GE [15] and PMR [6], they only provide methods for two modalities. Hence, to have a sufficient comparison, we retain the core uni-modal balancing strategy of OGM-GE and PMR, and extend them to more than two modality cases. Based on results in the left part of Table 3, our method is not limited by the number of modalities and remains effective in this case.

### 4.6    Comparison in scarcely informative modality case

As analyzed before, in existing imbalanced multimodal methods, the prevailing view is that the modality exhibiting superior predictive ability is considered the
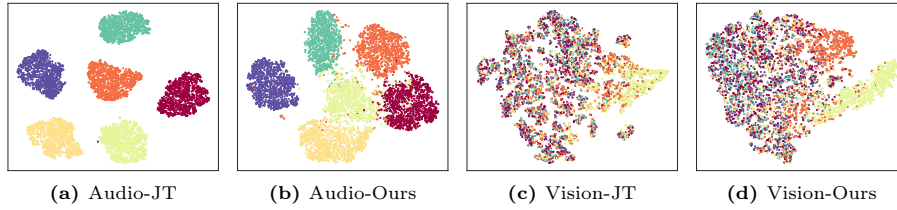
**(a)** Audio-JT          **(b)** Audio-Ours          **(c)** Vision-JT          **(d)** Vision-Ours

**Fig. 3:** Uni-modal representation visualization by t-SNE [12] on CREMA-D dataset. The categories are indicated in different colors. JT denotes for Joint-training.

"well-learnt" modality, while the remaining "worse-learnt" modalities require additional training during uni-modal balancing. Therefore, the existing imbalanced multimodal methods will fail when facing the case that modalities naturally have scarcely label-related information and more noise. This kind of scarcely informative modality will be recognized as the "worse-learnt" during training, and these methods will explicitly enhance their learning. However, enhancing their training offers no extra advancement and may even prompt the model to memorize more noise, affecting its effectiveness.

To validate this problem, we consider the scarcely informative modality case. We modify the audio data of the CREMA-D dataset, adding extra white Gaussian noise to make it noisier and scarcely discriminative. Based on the results shown in Table 4, all these imbalanced multimodal methods suffer a decline in performance when compared to the joint-training baseline, even the G-Blending [20] method with additional uni-modal modules and learning objectives. But our method continues to secure significant enhancement in this challenging scarcely informative modality case. The reason could be that this scarcely informative modality is often with both low-quality training representation and validation representation, due to the existence of much irrelevant noise. It has a small purity gap. Then, its encoder will be re-initialized with a slight percentage with our method, which is helpful to avoid the over-learning of the noisy data, even potentially improving model generalization [1,3,31]. Hence our method can well handle the scarcely informative modality case. In addition, there are related cases where one modality is not scarcely informative, but it is still noticeably less informative than others. For example, in the UCF-101 dataset, the accuracy of the individually-trained optical flow model is 58.9, and the individually-trained RGB model is 73.2. Our method also maintains superior, as shown in Table 1.

### 4.7   Uni-modal representation quality analysis

Beyond the comparison in overall multimodal performance, we also evaluate the uni-modal representation quality of our method to comprehensively reflect how well the imbalanced multimodal learning method is addressed. In terms of quantitative analysis, we fine-tune a new uni-modal classifier for the trained uni-modal encoder. Results are shown in Figure 1b and Figure 1c. Different from
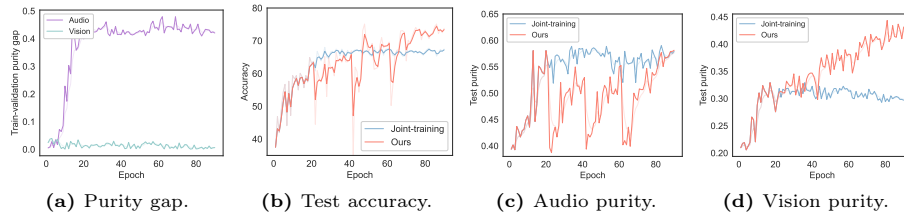
**(a)** Purity gap.     **(b)** Test accuracy.     **(c)** Audio purity.     **(d)** Vision purity.

**Fig. 4: (a):** The purity gap between training and validation representation. **(b):** Changes in test accuracy during training. **(c&d):** The purity of test representation. All results are based on the CREMA-D dataset.

other imbalanced multimodal learning methods that ignore or disrupt the training of well-learnt modalities, our method demonstrates an ideal enhancement for all modalities. In addition to the quantitative analysis, we also performed a qualitative analysis of uni-modal representation. As shown in Figure 3, we visualize the uni-modal representation by t-SNE [12] method, and have a comparison with the joint-training baseline. For the joint-training baseline, the audio modality is greatly separable, but the vision modality is with poor separability. In contrast, the audio representation separability of our method is ideal, although is slightly worse than the joint-training baseline. And the representation of vision modality has a noticeable improvement. The reason could be that our Diagnosing & Re-learning strategy can avoid the over-training of well-learnt modality while preserving its discriminative ability, and simultaneously encourage the training of other modalities. These quantitative and qualitative results demonstrate that our method effectively takes into account all modalities during balancing uni-modal learning.

### 4.8   Purity and accuracy analysis

In our diagnosing process, the learning state of each modality is estimated based on the separability discrepancy between training and validation representation space. And the separability is assessed by the purity of clustering results. In Figure 4a, we record the purity gap between the training and validation representation of the CREMA-D dataset. Based on the results, the audio modality is well-learnt with a huge purity gap while the vision modality is under-fitting. Therefore, during training, the audio encoder tends to be more greatly re-initialized. In experiments, we conduct the proposed method per 20 epochs for the CREMA-D dataset. The changes in test accuracy during training are shown in Figure 4b. Our method re-initializes the uni-modal encoders, especially one of well-learnt modality. This way firstly results in a sudden drop in performance due to model's reliance on the well-learnt modality, and the performance is progressively recovered and enhanced after re-learning the data. Besides the multimodal accuracy, we also observe the purity of uni-modal test representation. Based on Figure 4c and Figure 4d, for the well-learnt audio modality, the purity of our method
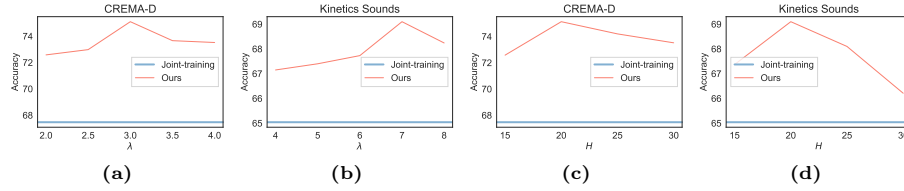
**Fig. 5:** Hyper-parameter sensitivity analysis of $\lambda$ in Equation 5 and Diagnosing & Re-learning frequency $H$ on CREMA-D and Kinetics Sounds datasets.

also has a similar trend to multimodal accuracy during training. The purity experiences a decreasing and increasing process. And for the under-fitting vision modality, its purity is continuously enhanced during training with our method, which indicates that our method effectively improves its learning.

### 4.9   Hyper-parameter sensitivity analysis

In this section, we conduct experiments to analyze two hyper-parameters $\lambda$ and $H$ in our method. Firstly, when assigning the re-initialization strength based on the uni-modal purity gap, the hyper-parameter $\lambda$ is introduced in Equation 5, to further control the re-initialization degree. Secondly, our Diagnosing & Re-learning strategy is conducted per $H$ epoch during training. Here we conduct experiments on both CREMA-D and Kinetics Sounds datasets about these two hyper-parameters. The results are shown in Figure 5. For $\lambda$, the results demonstrate that the Kinetics Sounds dataset tends to need a greater re-initialization degree than the CREMA-D dataset, but all these values consistently outperform the joint-training baseline. Also, for the re-initialization frequency $H$, performance is consistently enhanced across different frequencies, and its selection also does not require significant effort.

## 5   Conclusion

In this paper, we first analyze the limitations of existing imbalanced multimodal learning methods. They ignore the intrinsic limitation of modality capacity and the training of well-learnt modality during balancing. These limitations result in their failure in scarcely informative modality cases and may cause a decrease in the representation quality of well-learnt modality. To this end, we propose the Diagnosing & Re-learning method. It evaluates uni-modal learning state without any additional modules, and balances uni-modal training by softly re-initializing encoders, benefiting all modalities. Our method not only successfully overcomes the former limitations, but also exhibits its flexibility with diverse multimodal frameworks, well alleviating the imbalanced multimodal learning problem.

**Discussion.** Most current imbalanced multimodal learning methods focus on classification tasks. How to well estimate modality discrepancy and alleviate the imbalance in more types of tasks, *e.g.,* regression tasks, are still under-explored.

## Acknowledgements

## References

1. Alabdulmohsin, I., Maennel, H., Keysers, D.: The impact of reinitialization on generalization in convolutional neural networks. arXiv preprint arXiv:2109.00267 (2021)
2. Arandjelovic, R., Zisserman, A.: Look, listen and learn. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 609–617 (2017)
3. Ash, J., Adams, R.P.: On warm-starting neural network training. Advances in neural information processing systems **33**, 3884–3894 (2020)
4. Baltrušaitis, T., Ahuja, C., Morency, L.P.: Multimodal machine learning: A survey and taxonomy. IEEE transactions on pattern analysis and machine intelligence **41**(2), 423–443 (2018)
5. Cao, H., Cooper, D.G., Keutmann, M.K., Gur, R.C., Nenkova, A., Verma, R.: Crema-d: Crowd-sourced emotional multimodal actors dataset. IEEE transactions on affective computing **5**(4), 377–390 (2014)
6. Fan, Y., Xu, W., Wang, H., Wang, J., Guo, S.: Pmr: Prototypical modal rebalance for multimodal learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20029–20038 (2023)
7. Huang, Y., Lin, J., Zhou, C., Yang, H., Huang, L.: Modality competition: What makes joint training of multi-modal network fail in deep learning?(provably). arXiv preprint arXiv:2203.12221 (2022)
8. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
9. Li, H., Li, X., Hu, P., Lei, Y., Li, C., Zhou, Y.: Boosting multi-modal model performance with adaptive gradient modulation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 22214–22224 (2023)
10. Liang, P.P., Lyu, Y., Fan, X., Wu, Z., Cheng, Y., Wu, J., Chen, L., Wu, P., Lee, M.A., Zhu, Y., et al.: Multibench: Multiscale benchmarks for multimodal representation learning. arXiv preprint arXiv:2107.07502 (2021)
11. Liang, P.P., Zadeh, A., Morency, L.P.: Foundations and recent trends in multimodal machine learning: Principles, challenges, and open questions. arXiv preprint arXiv:2209.03430 (2022)
12. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(11) (2008)
13. MacQueen, J., et al.: Some methods for classification and analysis of multivariate observations. In: Proceedings of the fifth Berkeley symposium on mathematical statistics and probability. vol. 1, pp. 281–297. Oakland, CA, USA (1967)
14. Nagrani, A., Yang, S., Arnab, A., Jansen, A., Schmid, C., Sun, C.: Attention bottlenecks for multimodal fusion. Advances in Neural Information Processing Systems **34**, 14200–14213 (2021)
15. Peng, X., Wei, Y., Deng, A., Wang, D., Hu, D.: Balanced multimodal learning via on-the-fly gradient modulation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8238–8247 (2022)

16. Qiao, S., Lin, Z., Zhang, J., Yuille, A.L.: Neural rejuvenation: Improving deep network training by enhancing computational resource utilization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 61–71 (2019)
17. Sehwag, V., Chiang, M., Mittal, P.: On separability of self-supervised representations. In: ICML workshop on Uncertainty and Robustness in Deep Learning (UDL). vol. 3 (2020)
18. Sokar, G., Agarwal, R., Castro, P.S., Evci, U.: The dormant neuron phenomenon in deep reinforcement learning. In: Proceedings of the 40th International Conference on Machine Learning (2023)
19. Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human actions classes from videos in the wild. arXiv preprint arXiv:1212.0402 (2012)
20. Wang, W., Tran, D., Feiszli, M.: What makes training multi-modal classification networks hard? In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 12695–12705 (2020)
21. Wei, Y., Feng, R., Wang, Z., Hu, D.: Enhancing multimodal cooperation via sample-level modality valuation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 27338–27347 (2024)
22. Wei, Y., Hu, D.: Mmpareto: boosting multimodal learning with innocent unimodal assistance. In: International Conference on Machine Learning (2024)
23. Wei, Y., Hu, D., Tian, Y., Li, X.: Learning in audio-visual context: A review, analysis, and new perspective. arXiv preprint arXiv:2208.09579 (2022)
24. Wong, K.C.: A short survey on data clustering algorithms. In: 2015 Second international conference on soft computing and machine intelligence (ISCMI). pp. 64–68. IEEE (2015)
25. Wu, N., Jastrzebski, S., Cho, K., Geras, K.J.: Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks. In: International Conference on Machine Learning. pp. 24043–24055. PMLR (2022)
26. Xu, P., Zhu, X., Clifton, D.A.: Multimodal learning with transformers: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence (2023)
27. Yadav, S.K., Tiwari, K., Pandey, H.M., Akbar, S.A.: A review of multimodal human activity recognition with special emphasis on classification, applications, challenges and future directions. Knowledge-Based Systems **223**, 106970 (2021)
28. Yang, Z., Wei, Y., Liang, C., Hu, D.: Quantifying and enhancing multi-modal robustness with modality preference. In: The Twelfth International Conference on Learning Representations (2024)
29. Ying, X.: An overview of overfitting and its solutions. In: Journal of physics: Conference series. vol. 1168, p. 022022. IOP Publishing (2019)
30. Zadeh, A., Zellers, R., Pincus, E., Morency, L.P.: Mosi: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos. arXiv preprint arXiv:1606.06259 (2016)
31. Zaidi, S., Berariu, T., Kim, H., Bornschein, J., Clopath, C., Teh, Y.W., Pascanu, R.: When does re-initialization work? In: Proceedings on. pp. 12–26. PMLR (2023)
32. Zhu, H., Luo, M.D., Wang, R., Zheng, A.H., He, R.: Deep audio-visual learning: A survey. International Journal of Automation and Computing **18**(3), 351–376 (2021)