# Asymmetric Reinforcing against Multi-modal Representation Bias

**Xiyuan Gao**[1,2], **Bing Cao**[1,2*], **Pengfei Zhu**[1], **Nannan Wang**[2], **Qinghua Hu**[1]

[1]College of Intelligence and Computing, Tianjin University, Tianjin, 300000, China
[2]The State Key Laboratory of Integrated Services Networks, Xidian University, Xi'an, 710000, China
{gaoxiyuan, caobing, zhupengfei, huqinghua}@tju.edu.cn, nnwang@xidian.edu.cn

## Abstract

The strength of multimodal learning lies in its ability to integrate information from various sources, providing rich and comprehensive insights. However, in real-world scenarios, multi-modal systems often face the challenge of dynamic modality contributions, the dominance of different modalities may change with the environments, leading to suboptimal performance in multimodal learning. Current methods mainly enhance weak modalities to balance multimodal representation bias, which inevitably optimizes from a partial-modality perspective, easily leading to performance descending for dominant modalities. To address this problem, we propose an **A**symmetric **R**einforcing method against **M**ulti-modal representation bias (**ARM**). Our ARM dynamically reinforces the weak modalities while maintaining the ability to represent dominant modalities through conditional mutual information. Moreover, we provide an in-depth analysis that optimizing certain modalities could cause information loss and prevent leveraging the full advantages of multimodal data. By exploring the dominance and narrowing the contribution gaps between modalities, we have significantly improved the performance of multimodal learning, making notable progress in mitigating imbalanced multimodal learning. Our code is available at https://github.com/Gao-xiyuan/ARM.

## Introduction

Multimodal learning has emerged as a pivotal area in the field of machine learning, leveraging data from multiple sources to enhance the performance of models. This approach has been particularly transformative in applications such as image and text analysis, speech recognition, and autonomous driving, where combining visual, auditory, and textual information leads to more robust systems and makes multimodal learning an exciting frontier with significant potential (Huang et al. 2021). Despite promising yields, multimodal learning faces a critical challenge: *imbalanced learning among different modalities*. In most scenarios, partial modalities, even a single modality, may dominate the learning process, leading to insufficient learning of other modalities. Some modalities may become hard to learn due to environmental interference or limited information, leading to a multimodal bias for easier-to-learn modalities (Wu et al.
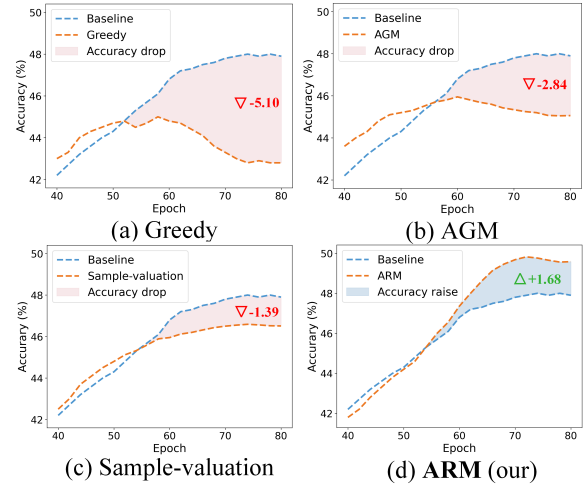
Figure 1: Accuracy curve of dominant modality compared with joint training baseline of imbalanced multimodal learning methods on Kinetics Sounds dataset. Other methods: Greedy (Wu et al. 2022), AGM(Li et al. 2023a), Sample-valuation (Wei et al. 2024).

2022), and multimodal learning may degrade to unimodal learning (Huang et al. 2022).

Imbalanced learning undermines the core objective of multimodal learning, which is to harness the complementary strengths of different data formats to achieve superior performance. In recent years, many extraordinary methods have been proposed to solve this problem, including canonical correlation analysis (Sun et al. 2020), random forest (Bi et al. 2020) and ensemble learning (Livne et al. 2018). Coupled with continuously optimized large-scale datasets and algorithm innovation, deep learning methods have shown significant promise in addressing modality imbalance (Lee, Lee, and Kim 2022; Das et al. 2023). The researchers attempted to balance the multimodal learning process through methods such as gradient modulation (Peng et al. 2022; Li et al. 2023a), collaborative learning (Rahate et al. 2022), and evaluation of modality contributions (Wei et al. 2024). However, these methods alleviate the imbalance by improving the representation of weak modalities from the uni-modal per-

spective alone, ignoring the connection between modalities, and not effectively utilizing all modalities. Although some methods (Zhang et al. 2024; Hu, Li, and Zhou 2022) consider cross-modal learning, they approach it from late fusion or modality preservation, without fully exploring the interrelationships between modalities, which limits their potential to improve model performance. Therefore, how to balance multimodal cooperation from a multimodal perspective remains an open question. Specifically, it is still expected to be addressed to narrow the contribution gaps between modalities and enhance the joint contribution of all modalities by exploring the interaction information between them.

To this concern, we have introduced a comprehensive valuation metric to evaluate the marginal contribution of each modality and the joint contribution of all modalities during learning for each sample. Mutual information (MI) originates from information theory used to measure the correlation between two random variables (Cover 1999). It represents the amount of information one variable contains about another and copes with capturing arbitrary dependency relationships, including linear, nonlinear, and higher-order relationships. This inspires us to use MI to measure the contribution of each modality to the learning process. To fully explore interaction information between modalities, we further utilize Conditional Mutual Information (CMI) to measure the reduction in uncertainty brought by introducing additional modalities on top of a uni-modal, thereby balancing multimodal learning without modality forgetting. Based on this, we propose an asymmetric enhancement method to dynamically alleviate imbalanced multimodal learning while maintaining the performance of dominant modalities. As shown in Fig. 1, most imbalanced multimodal learning methods exhibit dominant modality-forgetting during the training process because their optimization does not pay sufficient attention to dominant modalities, failing to maintain performance on these modalities. In contrast, based on the interrelationships between modalities, our method not only reasonably reduces the contribution disparity between modalities but also enhances the performance of each modality, overcoming the modality-forgetting. The main highlights of our study are as follows:

- We propose a mutual information-based valuation metric (MIV) to measure the marginal contribution of each modality and the joint contribution of all modalities in a sample with interrelation between modalities.

- Based on MIV, we propose an asymmetric reinforcement framework for multimodal representation bias, which dynamically narrows the contribution gaps between modalities. By continuously focusing on the dynamically changing dominance of different modalities, we mitigate modality forgetting and enhance the overall performance.

- We first reveal modality contributions from a multimodal perspective, each modality makes a positive and unique contribution to the multimodal systems. Extensive experiments validated our superiority on various multimodal classification datasets against the SOTAs.

## Related Works

### Imbalanced Multimodal Cooperation

Most multimodal learning often struggles with modality bias, where the dominant modality overshadows the others, leading to suboptimal performance. Recent advancements have focused on addressing this phenomenon through prototypical network (Fan et al. 2023), gradient modulation (Fu et al. 2023; Peng et al. 2022), and distilling knowledge (Pan et al. 2024; Du et al. 2021), dynamically weighing the importance of each modality based on task relevance or transferring knowledge from well-trained models, helping to mitigate imbalance. Evaluation methods (Koh et al. 2024; Yu et al. 2023), especially SHAPE (Hu, Li, and Zhou 2022) and Sample-valuation (Wei et al. 2024) novelly encourage balanced learning by improving the optimization of worse score modalities. Despite these advances, challenges remain in achieving truly balanced multimodal learning, most of these methods fall short by only enhancing weaker modalities without considering the intricate relationships between them. In contrast, we provide an asymmetric reinforcement strategy that dynamically alleviates multimodal bias based on contribution estimation without modality forgetting. This approach not only reduces the contribution disparity between modalities but also enhances overall multimodal cooperation, leading to improved performance across various multimodal classification datasets.

### Mutual Information in Machine Learning

Mutual Information (MI) has been a fundamental concept in information theory and its applications in machine learning (Haghifam et al. 2020; Hadizadeh et al. 2024), which highlights the dependency between variables. In machine learning, MI is widely used for feature selection and representation learning. Early techniques (Covert et al. 2023; Stutts et al. 2023) utilized MI to identify the most relevant features for predictive modeling, improve model performance by removing redundant or irrelevant features, and allow models to focus on the most informative features. In deep learning, MI has been instrumental in unsupervised learning and generative models (Larsson et al. 2019). Techniques like InfoGAN (Chen et al. 2016) leverage MI to improve the quality of generated samples and the robustness of models. Some variational autoencoders mutations (Pan, Long, and Pan 2023) use MI to learn a latent representation that captures the underlying structure of the data while ensuring independence between latent variables. Furthermore, MI neural estimation (Kim et al. 2022) has been introduced to efficiently estimate MI between high-dimensional variables, enabling more accurate learning in complex models. Recent advancements also include using MI in knowledge distillation (Chen et al. 2023) and domain adaptation (Wen et al. 2024), where understanding the information flow between different domains or causal variables is crucial. Overall, mutual information continues to be a powerful tool in enhancing the capabilities of machine learning models, driving us to use mutual information to measure modal benefits. To the best of our knowledge, we for the first time utilize mutual information to handle imbalanced multimodal learning.
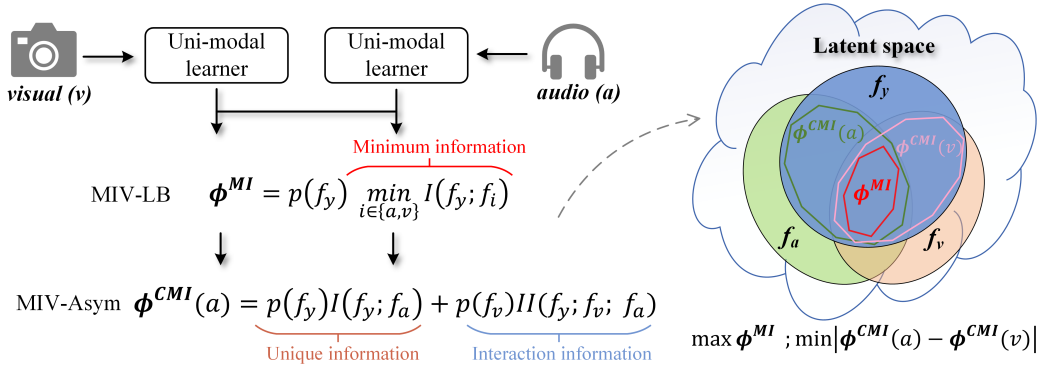
Figure 2: **Left:** The Lower Bound joint contribution (MIV-LB) of all modalities and the Asymmetric marginal contribution (MIV-Asym) of each modality are estimated by $\phi^{MI}$ and $\phi^{CMI}$, respectively, serving as the basis for asymmetric reinforcement. $f_{\mathcal{Y}}$ is feature-level fusion result, $p$ is the accurate production. **Right:** Representation of features in the latent space. We minimize the diversities in $\phi^{CMI}$ to balance multimodal learning while maximizing $\phi^{MI}$ to enhance multimodal performance.

## Methods

### Preliminary

In an interactive system, we can obtain partial information about one variable $X$ by observing another variable $Y$, thereby reducing the uncertainty of the former. The extent of this uncertainty reduction can be considered a measure of contribution and can be quantified using Mutual Information (MI). Using the basic relationship between the MI and entropy $H(\cdot)$ (Cover 1999), the algorithm for MI can be defined as the individual entropy of $X$, minus the conditional entropy of $X$ given $Y$. Following this approach, we can derive the formula for MI and Normalized MI (NMI) as follows:

$$I(X;Y) = \sum_{x \in X} \sum_{y \in Y} \mathcal{P}(x,y) \log \frac{\mathcal{P}(x,y)}{\mathcal{P}(x)\mathcal{P}(y)}, \quad (1)$$

$$NMI(X;Y) = \frac{I(X;Y)}{\sqrt{H(X)H(Y)}}. \quad (2)$$

Considering a multimodal classification task, a sample with $m$ modalities is represented as $\mathcal{X} = \{x^1, x^2, \ldots, x^m\}$, which can be regarded as a multimodal pair, and $y$ is the ground truth label of sample $\mathcal{X}$. Denote a uni-modal encoder as $\mathcal{E}(\cdot)$, the classification head as $\mathcal{H}(\cdot)$. The feature of the $i$-th modality extracted by the encoder is $f_{x^i} = \mathcal{E}(x^i)$. When taking $\mathcal{X}$ as the input for multimodal learning, the feature-level fusion output is $f_{\mathcal{Y}} = \cup f_{x^i}, x^i \in \mathcal{X}$, the final prediction is $\hat{y} = \mathcal{H}(f_{\mathcal{Y}})$. Notably, in this multimodal classification task, the features $f_{x^i}$ of the multimodal pair $\mathcal{X}$ are fused to obtain $f_{\mathcal{Y}}$. Subsequently, $f_{\mathcal{Y}}$ is used to make the final prediction $\hat{y}$, and the parameters of $\mathcal{E}(\cdot)$ are optimized by backpropagation based on $\hat{y}$, that is, $f_{\mathcal{Y}}$ further applied to $f_{x^i}$. Hence, a system characterized by the mutual interaction between $f_{x^i}$ and $f_{\mathcal{Y}}$ is constituted.

### Valuation Metric without Modality Forget

When the number of variables in an interactive system increases, such as in multimodal learning, where features $f_{\mathcal{X}} = \{f_{x^1}, f_{x^2}, \ldots, f_{x^m}\}$ from $m$ modalities jointly influence the fusion result $f_{\mathcal{Y}}$, using mutual information can become challenging. Inspired by exhaustively decomposing in a multivariate system (Williams and Beer 2010), even in cases where multiple source variables jointly influence a single variable, we can still compute the MI: $I(f_{\mathcal{Y}}; f_{x^i})$ for each $f_{x^i} \in f_{\mathcal{X}}$ with $f_{\mathcal{Y}}$ separately. Notably, Eq. (1) is non-negative, so it has a positive contribution to learning each modality. The MI between $f_{\mathcal{X}}$ and $f_{\mathcal{Y}}$ can be expressed as:

$$I(f_{\mathcal{Y}}; f_{\mathcal{X}}) = \sum_{\hat{y} \in f_{\mathcal{Y}}} \sum_{x \in f_{\mathcal{X}}} \mathcal{P}(x, \hat{y}) \log \frac{\mathcal{P}(x, \hat{y})}{\mathcal{P}(x)\mathcal{P}(\hat{y})}, \quad (3)$$

through observing $f_{\mathcal{Y}}$, the distribution of $f_{\mathcal{X}}$ changes from $\mathcal{P}(x)$ to $\mathcal{P}(x|\hat{y})$, we can capture the knowledge of $f_{\mathcal{X}}$ after the observation, the positive contribution in Eq. (3) is where predicting the ground truth label $y$, that is:

$$I(f_{\mathcal{Y}} = y; f_{\mathcal{X}}) = \sum_{\boldsymbol{x} \in f_{\mathcal{X}}} \mathcal{P}(\boldsymbol{x}|y) \log \frac{\mathcal{P}(y \mid \boldsymbol{x})}{\mathcal{P}(y)}. \quad (4)$$

**Theorem 1.** *In multimodal learning with m modalities, each modality can provide a **positive** and **unique** contribution to accurate prediction. i.e., $I(f_{\mathcal{Y}} = y; f_{x^i}) \neq I(f_{\mathcal{Y}} = y; f_{x^j})$, for any $x^i, x^j \in \mathcal{X}, i \neq j$. Naturally, neglecting the learning of any modality will result in information loss.* (The specific theoretical proof process is provided in the Appendix.)

Based **Theorem 1**, we propose a valuation metric to measure the marginal contribution of each modality in a sample $\mathcal{X}$, i.e. $\phi(x^i)$ and further derive the joint contribution of all modalities in that sample, i.e. $\phi(\mathcal{X})$. This serves as the foundation for asymmetric enhancement.

**Lower bound of joint contribution $\phi(\mathcal{X})$.** NMI between uni-modal and fused feature $NMI(f_{\mathcal{Y}}; f_{x^i})$ can be understood as the expected contribution value of all possible predictions from $f_{\mathcal{Y}}$ when $f_{x^i}$ is given, and it can be expressed as Eq. (5). For clarity, we use $I$ to represent $NMI$.

$$I(f_{\mathcal{Y}}; f_{x^i}) = \sum_{\hat{y}}^{N} p(f_{\mathcal{Y}} \to \hat{y}) I(f_{\mathcal{Y}}; f_{x^i}), \quad (5)$$

where $p(f_{\mathcal{Y}} \to \hat{y})$ represent the probability that $f_{\mathcal{Y}}$ makes the final prediction of class $\hat{y}$, $N$ is the number of categories. As we adopt *Softmax*, $\sum_{\hat{y}}^{N} p(f_{\mathcal{Y}} \to \hat{y}) = 1$. Therefore, based on the MI, the contribution of the model's accurate prediction provided by $i$-th modality can be written as:

$$\phi^{MI}\left(x^i\right) = p\left(f_{\mathcal{Y}} \to y\right) I\left(f_{\mathcal{Y}}; f_{x^i}\right). \quad (6)$$

Similarly, observing $j$-th modality ($j \neq i$) can also contribute to an extent that $f_{\mathcal{Y}}$ makes the accurate prediction $y$. Hence, the lower bound of joint contribution for all modalities in sample $\mathcal{X}$ is:

$$\phi^{MI}(\mathcal{X}) = p\left(f_{\mathcal{Y}} \to y\right) \min_{i \in \{1, \dots m\}} I\left(f_{\mathcal{Y}}; f_{x^i}\right). \quad (7)$$

It represents the minimum contribution value that each modality can provide for the model's accurate prediction. $\phi^{MI}$ has several properties: Firstly, its value range is $[0,1]$. Secondly, $\phi^{MI}$ is less than or equal to $I(f_{\mathcal{Y}}; f_{x^i})$ for all $i \leq m$. Finally, in the training phase, by incorporating $\phi^{MI}$ into the loss function and using gradient descent to maximize $\phi^{MI}$, thus each iteration moves towards increasing mutual information, ensuring the convergence of the lower bound.

**Estimating marginal contribution $\phi(x^i)$.** Although we defined the lower bound joint contribution of sample $\mathcal{X}$, the interrelationships between modalities are ignored, which prevents us from fully leveraging the advantages of multi-modal learning. As one would hope, given the presence of variable $Z$, the impact of introducing an additional variable $Y$ on $X$ can be measured using Conditional Mutual Information (CMI). The formulas for CMI and Normalized CMI (NCMI) are as follows:

$$CMI(X;Y|Z) = \sum_{x \in X} \sum_{y \in Y} \sum_{z \in Z} \mathcal{P}(x,y,z) \log \frac{\mathcal{P}(x,y|z)}{\mathcal{P}(x|z)\mathcal{P}(y|z)},$$
$$= \mathbb{E}_Z D_{KL}\left[\mathcal{P}(x,y|z) \| \mathcal{P}(x|z)\mathcal{P}(y|z)\right] \quad (8)$$
$$NCMI(X;Y|Z) = \frac{CMI(X;Y|Z)}{\sqrt{H(X|Z)H(Y|Z)}}. \quad (9)$$

In a complete modality set $\mathcal{X}$, when we choose the $x^i$ and $x^j$ to calculate MI with the fusion result separately, the mutual information of $f_{x^i}$ is $I(f_{\mathcal{Y}}; f_{x^i})$, and the conditional mutual information of $f_{x^j}$ given $f_{x^i}$ is $I(f_{\mathcal{Y}}; f_{x^j}|f_{x^i})$, vice versa. Consequently, the change in contribution value that modality $x^j$ causes to modality $x^i$ is the Interaction Information (II):

$$II(f_{\mathcal{Y}}; f_{x^j}; f_{x^i}) = I(f_{\mathcal{Y}}; f_{x^j}) - NCMI(f_{\mathcal{Y}}; f_{x^j}|f_{x^i}). \quad (10)$$

With Eq. (10), we can estimate the marginal contribution of $x^i$ based on considering all modalities as follows:

$$\phi^{CMI}(x^i) = p\left(f_{\mathcal{Y}} \to y\right) I\left(f_{\mathcal{Y}}; f_{x^i}\right)$$
$$+ \sum_{j \neq i}^{m} p(f_{x^j} \to y) II(f_{\mathcal{Y}}; f_{x^j}; f_{x^i}), \quad (11)$$

where $p(f_{x^j} \to y)$ can be regarded as a dynamic modality-specific weight of $j$-th modality, which can heighten the model's robustness (Yang et al. 2024). Furthermore, the joint contribution of the complete modality set from sample $\mathcal{X}$ can be expressed as:

$$\phi^{CMI}(\mathcal{X}) = \frac{1}{m} \sum_{i=1}^{m} \phi^{CMI}(x^i). \quad (12)$$

$\phi^{CMI}$ has several advantages: Firstly, it considers the impact of each modality from sample $\mathcal{X}$, ensuring that there is no modality omission during learning. Secondly, its value range is $[0, m]$, allowing it to be flexibly incorporated into loss functions or regularization as an optimization technique. Finally, averaging reasonably reflects the salient characteristics of the overall modalities, preventing the landslide victory of certain modalities while suppressing the occurrence of outliers.

**Asymmetric Reinforcement Strategies**

**Dynamic Feature-level Fusion.** Considering real-world factors, due to the primacy effect, the effect of the first term in Eq. (11) will be amplified. In other words, $\phi^{CMI}(x^i)$ reflects the importance to accurate prediction of $i$-th modality. We can use this as the specific-modal fusion weight during the fusion phase. Generally, higher $\phi^{CMI}(x^i)$ values represent more positive impacts on the model, thus the Fusion Weight (FW) of $i$-th modality can be denoted as:

$$FW^i = \frac{\phi^{CMI}(x^i)}{\phi^{CMI}(\mathcal{X})}, \quad (13)$$

where $FW^i$ works during the training phase and will take effect in the next epoch.

**Balanced Min-Max Loss.** Examining the expression of Eq. (7), (12), it is evident that $\phi^{MI}$ and $\phi^{CMI}$ represent the minimum contribution and comprehensive contribution that complete modalities for the model's accurate prediction, respectively. For the former, maximizing $\phi^{MI}$ enables the model to learn the most beneficial aspect of each modality for accurate prediction, and for the latter, we can use the Mean Absolute Error (MAE) to minimize $MAE(\phi^{CMI})$, thereby narrowing the marginal contribution gap between modalities.

$$\mathcal{L}_{\phi^{MI}} = 1 - \phi^{MI}(\mathcal{X}), \quad (14)$$
$$\mathcal{L}_{\phi^{CMI}} = \frac{\sum_{i=1}^{m} \left|\phi^{CMI}(x^i) - \phi^{CMI}(\mathcal{X})\right|}{\phi^{CMI}(\mathcal{X})}. \quad (15)$$

It should be noted that Eq. (14) cannot directly participate in the gradient backward process of gradient descent optimization since the min function is not globally differentiable. To this end, we use smooth approximation (Nielsen and Sun 2016) to make it differentiable:

$$\min_{i \in \{1, \dots m\}} I^i = \max_{i \in \{1, \dots m\}} (-I^i) \approx \log\left(\sum_{i=1}^{m} e^{-I^i}\right). \quad (16)$$

The overall loss function of ARM is formulated as Eq. (17), where $\mathcal{L}_{CE}$ denotes the cross-entropy loss, $\lambda_1$ and $\lambda_2$ are trade-off parameters.

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_1 \mathcal{L}_{\phi^{MI}} + \lambda_2 \mathcal{L}_{\phi^{CMI}}. \quad (17)$$

**Dynamic Sample-level Re-sample.** Following the analysis in **Theorem 1** and specific theoretical in (Wei et al. 2024), enhancing the discriminative ability of lower-contribution modality can expand its contribution. We propose to resample all modalities of lower joint contribution sample $\mathcal{X}$ more frequently during training. After each modality valuation by MIV, we can dynamically determine the re-sampling frequency with $\phi^{CMI}$ to enhance contribution, where re-sample frequency of sample $\mathcal{X}$ is:

$$s(\mathcal{X}) = \mathcal{F}_s(\phi^{CMI}(\mathcal{X})), \tag{18}$$

where $\mathcal{F}_s$ is a monotonically decreasing function, the lower-contribution sample $\mathcal{X}$ is re-trained with a resample frequency inversely proportional to its joint contribution. It is worth noting that different from (Wei et al. 2024), our re-sampling strategy is from a multimodal perspective, which dynamically adjusts the sampling frequency of all modalities in $\mathcal{X}$. This ensures that no information is lost during training, while the loss function $\mathcal{L}_{\phi^{MI}}$ guarantees targeted learning for lower-contribution modalities.

## Experiments

### Datasets and Implementation Details

**Kinetic Sounds (KS)** (Arandjelovic and Zisserman 2017) is a specifically designed action recognition dataset for research in audio-visual learning, particularly focusing on the relationship between actions and corresponding sounds. KS is composed of YouTube videos; all videos are cropped to within 10 seconds around the action. KS includes approximately 23k video clips with 31 categories.
**UCF-51** is a subset of UCF-101 (Soomro, Zamir, and Shah 2012) with two modalities, RGB and optical flow, containing 6,845 video clips across 51 diverse action categories. Mostly sourced from YouTube, it features varying conditions such as different camera angles and lighting, making it challenging and realistic for real-world applications.
**UPMC Food-101** (Wang et al. 2015) is a comprehensive dataset for food recognition, consisting of 101,000 images accompanied by corresponding texts across 101 food categories. Each category includes 750 images for training and 250 images for testing.

**Implementation Details.** Unless otherwise specified, ResNet-18 is used as the backbone in the experiments and trained from scratch. Encoders used for UCF-51 are ImageNet pre-trained. For Food-101, a ViT-based model is used as the vision encoder, and a BERT-based model is used as the text encoder by the pre-trained. Before modality valuation, a warm-up stage is employed for all experiments. During training, we use Stochastic Gradient Descent (SGD) with a batch size of 64. We set the initial learning rate, weight decay, and momentum parameters to $10^{-3}$, $5 \times 10^{-4}$, and 0.9, respectively. The experiments are conducted on Huawei Atlas 800 Training Server with CANN and NVIDIA 4090 GPU. More details of implementation and experiment analysis are provided in the Appendix.

| Model | KS (Acc.) | UCF-51 (Acc.) |
|---|---|---|
| Concatenation | 59.61 | 68.23 |
| Summation | 59.53 | 67.62 |
| OGM-GE (CVPR 2022) | 60.70 | 71.66 |
| Greedy (ICML 2022) | 59.86 | 71.53 |
| QMF (ICML 2023) | 63.78 | 73.48 |
| PMR (CVPR 2023) | 63.86 | 74.80 |
| Sample-val. (CVPR 2024) | 65.33 | 75.12 |
| Modality-val. (CVPR 2024) | 65.10 | 74.39 |
| MLA (CVPR 2024) | 65.21 | 76.01 |
| **ARM** | 66.52 | 75.60 |

Table 1: Accuracy of imbalanced multimodal learning methods, where red and blue indicate the best/runner-up performance. Results are reported in percentage ($\%$).
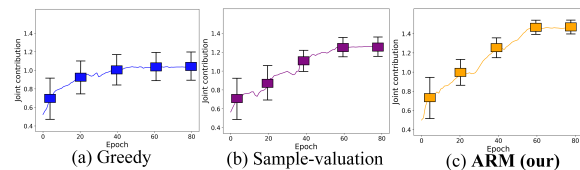


Figure 3: Comparison of the narrowing trend of uni-modality contribution gaps on the UCF-51 dataset.

### Comparison with Imbalanced Multimodal Learning Methods

In this section, we compared ARM with advanced imbalanced multimodal learning methods to answer **Q1:** *How does ARM narrow the modality contribution gap?*

Fig. 3 illustrates the trend of narrowing contribution gaps across different methods. Traditional methods, like Greedy, have shown limited improvement in closing the contribution disparity between modalities, with only a slight narrowing in the contribution gap as training progresses. Sample-valuation demonstrates more consistent shrink, yet the gap remains noticeable across epochs. In contrast, ARM achieves a marked and consistent reduction in modality contribution gaps, indicating a more balanced learning process. This consistent improvement shows ARM's ability to maintain equitable contribution from all modalities, which is crucial for robust multimodal learning.

Table 1 further reinforces this conclusion. ARM consistently outperforms other state-of-the-art methods, i.e., Greedy (Wu et al. 2022), OGM-GE (Peng et al. 2022), QMF (Zhang et al. 2023), PMR (Fan et al. 2023), Sample-valuation, Modality-valuation (Wei et al. 2024), and MLA (Zhang et al. 2024), achieving the competitive accuracy scores of 66.52% and 75.60%, respectively. Other approaches, like QMF and PMR, show decent performance but still fall short in balancing modality contributions, leading to suboptimal accuracy. Due to the different design focus, MLA performs better in handling temporal optical flow data in the UCF-51 dataset. Sample-valuation achieves competitive results but cannot match the balance achieved by ARM, which is evident from the joint contribution trends shown in Fig. 3. The advantage of ARM lies in its dual focus:

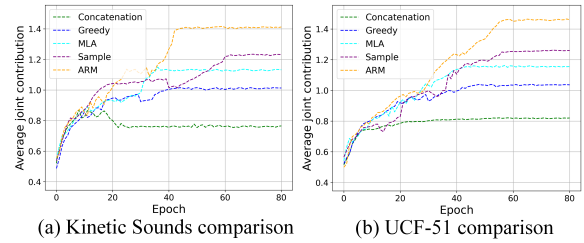| Model | KS (Acc.) | Food-101 (Acc.) |
|---|---|---|
| Concatenation | 59.61 | 82.38 |
| Summation | 59.53 | 82.63 |
| Decision fusion | 60.12 | 83.71 |
| FiLM (AAAI 2018) | 59.33 | 82.34 |
| BiGated (AAAI 2018) | 60.79 | 86.71 |
| Dynamic Fusion (CVPR 2023) | 63.21 | 90.83 |
| PMF (ICCV 2023) | 64.33 | 91.56 |
| TransFusion (ICLR 2024) | 65.40 | 91.22 |
| ARM | 66.52 | 93.36 |

Table 2: Comparison with multimodal fusion methods.



Figure 4: Average joint contribution of all modalities overall training samples during training for Greedy, MLA, Sample-valuation, and our ARM on the KS and UCF-51 datasets.

minimizing modality imbalances while maximizing overall performance. By effectively narrowing the contribution gaps between modalities, ARM prevents any single modality from dominating or being neglected, leading to a more cohesive and effective multimodal representation.

## Comparison with Multimodal Fusion Methods

Table 2 compares the performance of various multimodal fusion methods on two datasets to answer **Q2:** *Can the proposed modules (e.g., dynamic feature-level fusion) effectively improve performance?*

Concatenation and Summation are baseline methods that simply merge the feature vectors, yielding moderate performance. More advanced techniques such as FiLM (Perez et al. 2018) and BiGated (Kiela et al. 2018) introduce interaction between modalities through modulation or gating mechanisms, resulting in eligible accuracy compared with the baseline. Dynamic Fusion (Xue and Marculescu 2023) incorporates adaptive fusion strategies, which inspire our dynamic feature-level fusion, adjusting how the modalities are combined during inference, which leads to substantial improvements, especially on the Food-101 dataset.

Among the recently proposed methods, PMF (Li et al. 2023b) and TransFusion (Imfeld et al. 2023) demonstrate the power of more sophisticated fusion techniques. PMF achieves strong performance by effectively managing modality-specific features. TransFusion, a transformer-based model, further refines this by better capturing the complex interactions between modalities, achieving runner-up results. Our ARM outperforms all other models on both datasets, achieving 66.52% accuracy on KS and 93.36% on Food-101, which is a significant improvement, particularly evident on the KS dataset, where it exceeds the runner-up by over 1 percentage point. ARM's success is attributed to its advanced asymmetric reinforcement strategy, which dynamically balances the learning from each modality, preventing the model from being biased toward the dominant modality. This ensures that both audio and visual information are utilized effectively, leading to superior performance in challenging multimodal tasks. Compared with the competing methods, ARM's ability to maintain high accuracy across different datasets demonstrates its robustness and adaptability, making it a standout choice for multimodal fusion tasks.

## Analysis of Modality Forget & Multimodal Cooperation

We report the results of a single modality and a combination of all modalities and further display the improvement of multimodal cooperation to answer **Q3:** *Compared to prior multimodal learning approaches, can ARM overcome modality forget and optimize multimodal cooperation?*

**Modality Forget.** Table 3 compares the performance of various models across multiple datasets, highlighting results for different modalities and their multimodal cooperation. Several models in the comparison exhibit a notable modality forget phenomenon, where optimizing for one weaker modality leads to a decrease in the performance of the dominant modality and achieves suboptimal results in the overall multimodal performance. For instance, on the KS dataset, models like BiGated and PMF show significant drops in performance for the visual modality compared to the audio, which in turn negatively impacts their multimodal accuracy. This trend is also observed on UCF-51, where models fail to balance the learning of RGB and optical flow modalities, leading to lower overall performance. The Sample-valuation model also shows a considerable drop across both visual and textual modalities on Food-101, which further highlights the issue of modality forgetting. Our proposed ARM consistently outperforms other models across all datasets, achieving the highest accuracy in both single and multimodal scenarios. Notably, ARM excels in preventing modality forgetting, as demonstrated by its superior performance across different modalities and their combinations.

**Multimodal Cooperation.** Fig. 4 illustrates the progression of multimodal joint contributions over epochs compared to Concatenation baseline for different models. The performance of the other methods indicates a relatively slower and less stable increase in the multimodal joint contribution over time. Greedy demonstrates some improvement but plateaus early, indicating that it struggles to maintain steady enhancement of multimodal cooperation. Sample and MLA show better performance than Concatenation and Greedy but still fall short compared to ARM, as they are unable to fully exploit the joint potential of multimodal learning. ARM exhibits a consistent and substantial increase in the multimodal average contribution, the chart shows that ARM not only achieves a higher overall contribution but also

| Dataset | | Conact. | Sum | BiGated | PMF | QMF | Sample | MLA | **ARM** |
|---|---|---|---|---|---|---|---|---|---|
| KS | (⋆) Audio | 47.35 | 46.21 | 44.11 (↓) | 45.82 (↓) | 47.56 | 46.02 (↓) | 49.20 | 49.95 |
| | Video | 23.65 | 22.78 | 22.08 (↓) | 25.65 | 36.82 | 42.67 | 41.30 | 44.86 |
| | Mutli | 59.61 | 59.53 | 60.79 | 64.33 | 63.78 | 65.33 | 65.21 | 66.52 |
| UCF-51 | (⋆) RGB | 60.13 | 59.80 | 57.39 (↓) | 58.13 (↓) | 56.20 (↓) | 57.01 (↓) | 64.81 | 63.29 |
| | OF | 29.62 | 28.81 | 25.67 (↓) | 36.21 | 40.51 | 42.33 | 41.26 | 43.19 |
| | Mutli | 68.23 | 67.62 | 70.87 | 72.09 | 73.48 | 75.12 | 76.01 | 75.60 |
| Food-101 | Image | 30.85 | 31.66 | 48.87 | 59.21 | 66.39 | 73.49 | 71.58 | 72.36 |
| | (⋆) Text | 81.68 | 80.84 | 78.51 (↓) | 79.66 (↓) | 82.10 | 84.43 | 86.42 | 86.86 |
| | Mutli | 82.38 | 82.63 | 86.71 | 91.56 | 91.67 | 90.85 | 93.31 | 93.36 |

Table 3: Comparison results on audio-video, RGB-optical flow, and image-text datasets. The performance of a single modality and the results of combining all modalities ("multiple") are listed. ⋆ denotes the dominant modality and ↓ indicates a performance drop compared with Concatenation or Sum baseline.

| $\mathcal{L}_{CE}$ | $\mathcal{L}_{\phi^{MI}}$ | $\mathcal{L}_{\phi^{CMI}}$ | KS | UCF-51 | Food-101 |
|---|---|---|---|---|---|
| ✓ | | | 63.88 | 70.03 | 88.52 |
| ✓ | ✓ | | 64.20 | 73.56 | 91.25 |
| ✓ | | ✓ | 65.19 | 72.10 | 89.78 |
| ✓ | ✓ | ✓ | 66.52 | 75.60 | 93.36 |

Table 4: Ablation study of loss function.



(a) Kinetic Sounds  (b) UCF-51

Figure 5: Curve of Balanced Min-Max Loss: the values are obtained from 5 training processes with the same initiations.

demonstrates a stable and continuous growth trend, indicating its robustness in learning and integrating information from various modalities.

ARM's success can be attributed to its dynamic asymmetric reinforcement strategy, which effectively balances the learning contributions from each modality based on their importance. By dynamically adjusting the contribution of each modality based on their importance and interaction with others, ARM ensures that no single modality dominates at the expense of others and allows ARM to maximize the joint contribution of all modalities, leading to superior performance in multimodal learning.

### The Effectiveness of Loss Function

This section answers the question: **Q4:** *Does our proposed Balanced Min-Max loss progress as expected?*

Fig. 5 demonstrates the effectiveness of the proposed $\mathcal{L}_{\phi^{MI}}$ and $\mathcal{L}_{\phi^{CMI}}$ in improving overall multimodal performance and alleviating imbalanced learning between modalities, respectively. The loss curves for both the KS and UCF-51 datasets show that incorporating the Balanced Min-Max loss consistently improves the overall model performance by ensuring better multimodal cooperation, leading to faster convergence and lower loss values. Table 4 further validates these observations with an ablation study. When only the $\mathcal{L}_{\phi^{MI}}$ is added, there is a noticeable increase in accuracy compared to the baseline (row 1). Additionally, the inclusion of the $\mathcal{L}_{\phi^{CMI}}$ specifically addresses modality imbalance by reducing the learning disparity between modalities. This is particularly important in scenarios where dominant modalities may overshadow weaker ones. The combined use of both $\mathcal{L}_{\phi^{MI}}$ and $\mathcal{L}_{\phi^{CMI}}$ achieves the best performance, demonstrating that our approach not only enhances overall accuracy but also maintains balanced contributions across all modalities. This synergy between the two losses highlights the strength of our method in multimodal learning.

### Conclusion

In this paper, we introduce a valuation metric to evaluate the marginal contributions of different modalities and the joint contributions of all modalities in a sample with a theoretical analysis of mutual information. Based on this, an asymmetric enhancement method named ARM is proposed to improve imbalanced multimodal learning while preventing modality forgetting. This provides a potential approach for balancing multimodal learning in real-world applications. Besides, there are some further discussions.

**Universality of Mutual Information.** Our method calculates mutual information after feature extraction, and the data dimension is lower, so we can directly calculate the marginal distribution and joint distribution. But when processing continuous data, discretization or kernel density estimation methods are required. These methods are relatively complex to implement and may lead to different results.

**Natural Conflict in Multimodal.** Multimodal data may contain some inherent conflicts. For example, for an RGB-Infrared sample *person* in foggy environments, two modalities may make vastly different predictions. Although ARM copes with mitigating modality conflicts by reducing the impact of modalities with incorrect predictions, it does not fundamentally resolve such conflicts. Therefore, it is expected to consider this natural conflict in the future work.

## Acknowledgments

## References

Arandjelovic, R.; and Zisserman, A. 2017. Look, listen and learn. In *Proceedings of the IEEE international conference on computer vision*, 609–617.

Bi, X.-a.; Hu, X.; Wu, H.; and Wang, Y. 2020. Multimodal data analysis of Alzheimer's disease based on clustering evolutionary random forest. *IEEE Journal of Biomedical and Health Informatics*, 24: 2973–2983.

Chen, M.; Xing, L.; Wang, Y.; and Zhang, Y. 2023. Enhanced multimodal representation learning with cross-modal kd. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11766–11775.

Chen, X.; Duan, Y.; Houthooft, R.; Schulman, J.; Sutskever, I.; and Abbeel, P. 2016. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29.

Colombo, P.; Chapuis, E.; Labeau, M.; and Clavel, C. 2021. Improving multimodal fusion via mutual dependency maximisation. *arXiv preprint arXiv:2109.00922*.

Cover, T. M. 1999. *Elements of information theory*. John Wiley & Sons.

Covert, I. C.; Qiu, W.; Lu, M.; Kim, N. Y.; White, N. J.; and Lee, S.-I. 2023. Learning to maximize mutual information for dynamic feature selection. In *International Conference on Machine Learning*, 6424–6447. PMLR.

Das, A.; Das, S.; Sistu, G.; Horgan, J.; Bhattacharya, U.; Jones, E.; Glavin, M.; and Eising, C. 2023. Revisiting modality imbalance in multimodal pedestrian detection. In *2023 IEEE International Conference on Image Processing (ICIP)*, 1755–1759. IEEE.

Du, C.; Li, T.; Liu, Y.; Wen, Z.; Hua, T.; Wang, Y.; and Zhao, H. 2021. Improving multi-modal learning with uni-modal teachers. *arXiv preprint arXiv:2106.11059*.

Fan, Y.; Xu, W.; Wang, H.; Wang, J.; and Guo, S. 2023. Pmr: Prototypical modal rebalance for multimodal learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20029–20038.

Fu, J.; Gao, J.; Bao, B.-K.; and Xu, C. 2023. Multimodal imbalance-aware gradient modulation for weakly-supervised audio-visual video parsing. *IEEE Transactions on Circuits and Systems for Video Technology*.

Hadizadeh, H.; Yeganli, S. F.; Rashidi, B.; and Bajić, I. V. 2024. Mutual Information Analysis in Multimodal Learning Systems. *arXiv preprint arXiv:2405.12456*.

Haghifam, M.; Negrea, J.; Khisti, A.; Roy, D. M.; and Dziugaite, G. K. 2020. Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. *Advances in Neural Information Processing Systems*, 33: 9925–9935.

Han, W.; Chen, H.; and Poria, S. 2021. Improving multimodal fusion with hierarchical mutual information maximization for multimodal sentiment analysis. *arXiv preprint arXiv:2109.00412*.

Hu, P.; Li, X.; and Zhou, Y. 2022. Shape: An unified approach to evaluate the contribution and cooperation of individual modalities. *arXiv preprint arXiv:2205.00302*.

Huang, Y.; Du, C.; Xue, Z.; Chen, X.; Zhao, H.; and Huang, L. 2021. What makes multi-modal learning better than single (provably). *Advances in Neural Information Processing Systems*, 34: 10944–10956.

Huang, Y.; Lin, J.; Zhou, C.; Yang, H.; and Huang, L. 2022. Modality competition: What makes joint training of multimodal network fail in deep learning?(provably). In *International conference on machine learning*, 9226–9259. PMLR.

Imfeld, M.; Graldi, J.; Giordano, M.; Hofmann, T.; Anagnostidis, S.; and Singh, S. P. 2023. Transformer fusion with optimal transport. *arXiv preprint arXiv:2310.05719*.

Kiela, D.; Grave, E.; Joulin, A.; and Mikolov, T. 2018. Efficient large-scale multi-modal classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Kim, J.-H.; Kim, Y.; Lee, J.; Yoo, K. M.; and Lee, S.-W. 2022. Mutual information divergence: A unified metric for multimodal generative models. *Advances in Neural Information Processing Systems*, 35: 35072–35086.

Koh, J. Y.; Lo, R.; Jang, L.; Duvvur, V.; Lim, M. C.; Huang, P.-Y.; Neubig, G.; Zhou, S.; Salakhutdinov, R.; and Fried, D. 2024. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*.

Larsson, M.; Stenborg, E.; Toft, C.; Hammarstrand, L.; Sattler, T.; and Kahl, F. 2019. Fine-grained segmentation networks: Self-supervised segmentation for improved long-term visual localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 31–41.

Lee, H. K.; Lee, J.; and Kim, S. B. 2022. Boundary-focused generative adversarial networks for imbalanced and multimodal time series. *IEEE Transactions on Knowledge and Data Engineering*, 34: 4102–4118.

Li, H.; Li, X.; Hu, P.; Lei, Y.; Li, C.; and Zhou, Y. 2023a. Boosting Multi-modal Model Performance with Adaptive Gradient Modulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 22214–22224.

Li, Y.; Quan, R.; Zhu, L.; and Yang, Y. 2023b. Efficient multimodal fusion via interactive prompting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2604–2613.

Liao, R.; Moyer, D.; Cha, M.; Quigley, K.; Berkowitz, S.; Horng, S.; Golland, P.; and Wells, W. M. 2021. Multimodal representation learning via maximization of local mutual information. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th*

*International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, 273–283. Springer.

Livne, M.; Boldsen, J. K.; Mikkelsen, I. K.; Fiebach, J. B.; Sobesky, J.; and Mouridsen, K. 2018. Boosted tree model reforms multimodal magnetic resonance imaging infarct prediction in acute stroke. *Stroke*, 49: 912–918.

Nielsen, F.; and Sun, K. 2016. Guaranteed bounds on information-theoretic measures of univariate mixtures using piecewise log-sum-exp inequalities. *Entropy*, 18: 442.

Pan, W.; Long, F.; and Pan, J. 2023. ScInfoVAE: interpretable dimensional reduction of single cell transcription data with variational autoencoders and extended mutual information regularization. *BioData Mining*, 16: 17.

Pan, Y.; Jiang, J.; Jiang, K.; and Liu, X. 2024. Disentangled-Multimodal Privileged Knowledge Distillation for Depression Recognition with Incomplete Multimodal Data. In *ACM Multimedia*.

Peng, X.; Wei, Y.; Deng, A.; Wang, D.; and Hu, D. 2022. Balanced multimodal learning via on-the-fly gradient modulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 8238–8247.

Perez, E.; Strub, F.; De Vries, H.; Dumoulin, V.; and Courville, A. 2018. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Rahate, A.; Walambe, R.; Ramanna, S.; and Kotecha, K. 2022. Multimodal co-learning: Challenges, applications with datasets, recent advances and future directions. *Information Fusion*, 81: 203–239.

Soomro, K.; Zamir, A. R.; and Shah, M. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*.

Stutts, A. C.; Erricolo, D.; Ravi, S.; Tulabandhula, T.; and Trivedi, A. R. 2023. Mutual information-calibrated conformal feature fusion for uncertainty-aware multimodal 3d object detection at the edge. *arXiv preprint arXiv:2309.09593*.

Sun, Z.; Sarma, P.; Sethares, W.; and Liang, Y. 2020. Learning relationships between text, audio, and video via deep canonical correlation for multimodal language analysis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 8992–8999.

Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9.

Wang, X.; Kumar, D.; Thome, N.; Cord, M.; and Precioso, F. 2015. Recipe recognition with large multimodal food dataset. In *2015 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, 1–6. IEEE.

Wei, Y.; Feng, R.; Wang, Z.; and Hu, D. 2024. Enhancing multimodal cooperation via sample-level modality valuation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 27338–27347.

Wen, L.; Chen, S.; Xie, M.; Liu, C.; and Zheng, L. 2024. Training multi-source domain adaptation network by mutual information estimation and minimization. *Neural Networks*, 171: 353–361.

Williams, P. L.; and Beer, R. D. 2010. Nonnegative decomposition of multivariate information. *arXiv preprint arXiv:1004.2515*.

Wu, N.; Jastrzebski, S.; Cho, K.; and Geras, K. J. 2022. Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks. In *International Conference on Machine Learning*, 24043–24055. PMLR.

Xue, Z.; and Marculescu, R. 2023. Dynamic multimodal fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2575–2584.

Yang, Z.; Wei, Y.; Liang, C.; and Hu, D. 2024. Quantifying and Enhancing Multi-modal Robustness with Modality Preference. In *The Twelfth International Conference on Learning Representations*.

Yu, W.; Yang, Z.; Li, L.; Wang, J.; Lin, K.; Liu, Z.; Wang, X.; and Wang, L. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.

Zhang, Q.; Wu, H.; Zhang, C.; Hu, Q.; Fu, H.; Zhou, J. T.; and Peng, X. 2023. Provable Dynamic Fusion for Low-Quality Multimodal Data. In *International Conference on Machine Learning*.

Zhang, X.; Yoon, J.; Bansal, M.; and Yao, H. 2024. Multimodal representation learning by alternating unimodal adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 27456–27466.

# Appendix

In Appendix, we provide more details, proofs and experiments, encompassing the following:

- A theoretical proof of Theorem 1, with the help of mutual information, we have proved the indispensability of each modality. The proof is elaborated in *Appx. A*.

- More details of experiment implementation, including data processing, experiment setting and algorithm flow, as detailed in *Appx. B*.

- Additional experimental comparisons, including more ablated experiment on ARM, more discussion on dynamic sample-level re-sample strategy, more comparisons with other MI-based methods, and the improvement of ARM on other fusion methods, etc, as presented in *Appx. C*.

## A. Proof

**Theorem 1.** *In multimodal learning with m modalities, each modality can provide a **positive** and **unique** contribution to accurate prediction. i.e., $I(f_\mathcal{Y} = y; f_{x^i}) \neq I(f_\mathcal{Y} = y; f_{x^j})$, for any $x^i, x^j \in \mathcal{X}, i \neq j$. Naturally, neglecting the learning of any modality will result in information loss.*

*Proof.* Let us consider 3 finite non-empty feature sets of a sample $\mathcal{X}$ with $m$ modalities: $A = \{f_{x^1}, f_{x^2}, \ldots, f_{x^m}\}$, $B = \{f_{x^1}, f_{x^2}, \ldots, f_{x^n}\}$, and $C = \{f_{x^{n+1}}, f_{x^{n+2}}, \ldots, f_{x^m}\}$, $n < m, A = B \cup C$, jointly affecting $Y = f_\mathcal{Y}$. Under the condition of ensuring accurate prediction, that is, $Y = y$, we have Eq. (19).

Consequently, the positive contribution provided by set $A$ must be greater than that of set $B$. Neglecting the learning of any modality will result in the loss of positive information. In other words, each modality in a sample contains a unique positive impact. This prompts us to optimize imbalanced multimodal cooperation without abandoning any modality.

## B. More Details

### Implementation Details.

For the KS dataset, the network was trained for 80 epochs using the SGD optimizer with a momentum of 0.9, a learning rate of 0.001, and a weight decay of 0.0005. The batch size was set to 64. All videos in the KS dataset were resized to have a short edge of 256 pixels, and the sampling frequency was set to one frame per second. For the UCF-51 dataset, the encoders were initialized with ImageNet pretrained weights, and the network was trained with an initial learning rate of 0.0005 and a batch size of 16. For the Food-101 dataset, the AdamW optimizer was used with a learning rate of 0.0001, and the backbone was also pre-trained on ImageNet. A warm-up phase was employed for all experiments, with the warm-up duration set to 10 epochs. All other hyperparameters were set to their default values as defined in PyTorch.

$$
\begin{aligned}
& I(Y = y; \boldsymbol{A}) - I(Y = y; \boldsymbol{B}) \\
&= \sum_{\boldsymbol{a}} \mathcal{P}(\boldsymbol{a} \mid y) \log \frac{\mathcal{P}(y, \boldsymbol{a})}{\mathcal{P}(y)\mathcal{P}(\boldsymbol{a})} \\
&\quad - \sum_{\boldsymbol{b}} \mathcal{P}(\boldsymbol{b} \mid y) \log \frac{\mathcal{P}(y, \boldsymbol{b})}{\mathcal{P}(y)\mathcal{P}(\boldsymbol{b})} \\
&= \sum_{\boldsymbol{b}} \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{b}, \boldsymbol{c} \mid y) \log \frac{\mathcal{P}(y, \boldsymbol{b}, \boldsymbol{c})}{\mathcal{P}(y)\mathcal{P}(\boldsymbol{b}, \boldsymbol{c})} \\
&\quad - \sum_{\boldsymbol{b}} \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{b}, \boldsymbol{c} \mid y) \log \frac{\mathcal{P}(y, \boldsymbol{b})}{\mathcal{P}(y)\mathcal{P}(\boldsymbol{b})} \quad (19) \\
&= \sum_{\boldsymbol{b}} \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{b}, \boldsymbol{c} \mid y) \log \frac{\mathcal{P}(\boldsymbol{b})\mathcal{P}(\boldsymbol{b})\mathcal{P}(y \mid \boldsymbol{b})\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y)}{\mathcal{P}(\boldsymbol{c}, \boldsymbol{b})\mathcal{P}(y, \boldsymbol{b})} \\
&= \sum_{\boldsymbol{b}} \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{b}, \boldsymbol{c} \mid y) \log \frac{\mathcal{P}(\boldsymbol{b})\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y)}{\mathcal{P}(\boldsymbol{b}, \boldsymbol{c})} \\
&= \sum_{\boldsymbol{b}} \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{b})\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y) \log \frac{\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y)}{\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b})} \\
&= \sum_{\boldsymbol{b}} \mathcal{P}(\boldsymbol{b}) \sum_{\boldsymbol{c}} \mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y) \log \frac{\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y)}{\mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b})} \\
&= \mathbb{E}_B D_{KL} \left[ \mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}, y) \| \mathcal{P}(\boldsymbol{c} \mid \boldsymbol{b}) \right] \geq 0
\end{aligned}
$$

### Algorithm Details.

The whole training pipeline is provided in Algorithm 1. ARM outlines a contribution enhancement strategy designed to address imbalanced multimodal learning challenges. It begins by initializing $\mathcal{D}^{rs}$ as $\mathcal{D}$. For each sample of multimodal inputs $\mathcal{X} = \{x^1, x^2, \ldots, x^m\}$ in $\mathcal{D}^{rs}$, ARM evaluates the multimodal joint contribution first. After the warm-up period, the core of the algorithm activates, and the contribution scores guide a dynamic feature-level fusion process. Next, the algorithm calculates the re-sampling frequency $s(\mathcal{X})$ for each sample and updates $\mathcal{D}^{rs}$ based on these frequencies. The process then repeats loss $\mathcal{L}_{total}$ computation and parameter updates using the modified dataset. This approach enhances learning by continuously adjusting contributions from each modality based on their mutual interactions, promoting balanced learning in multimodal scenarios.

## C. More Discussions

### Analysis of Different Re-sample Frequency

In this section, we provide a comparison of the results with other re-sample methods and different sampling frequencies to answer **Q5:** *Does our dynamic sample-level re-sample strategy really work?*

We compare with two related re-sample settings, Random re-sample is to randomly re-sample input of each sample with the same frequency, Inverse re-sample is only re-sampling the sample with higher contribution. Our proposed dynamic sample-level re-sample (DSR), sampling function $\mathcal{F}_s = round < k\phi^{CMI}(\mathcal{X}) - km >$, where $\mathcal{X}$ represents

Algorithm 1: Asymmetric reinforcement strategies

**Require:** Original training dataset $\mathcal{D}$, training dataset with re-sample $\mathcal{D}^{rs}$, number of modalities $m$, loss fuction $\mathcal{L}_{total}$, model parameters $\theta$, training epoch $T$, warm-up epoch $F$.

1: **for** $t = 0, \ldots, T-1$ **do**
2:     Initialize $\mathcal{D}^{rs} = \mathcal{D}$;
3:     **for** each sample $\mathcal{X} = \{x^1, x^2, \ldots, x^m\}$ in $\mathcal{D}^{rs}$ **do**
4:         Valuate multimodal joint contribution $\phi^{MI}(\mathcal{X}), \phi^{CMI}(\mathcal{X})$ with Eq. (7), (12);
5:         **if** $t < F$ **then**
6:             Compute the loss $\mathcal{L}_{total}$ following Eq. (17);
7:             Update parameters $\theta$ with dataset $\mathcal{D}^{rs}$;
8:         **else**
9:             Dynamic feature fusion with Eq. (13);
10:           Get re-sample frequency $s(\mathcal{X})$ with Eq. (18);
11:           Add $\mathcal{X}$ with frequency $s(\mathcal{X})$ into $\mathcal{D}^{rs}$;
12:           Compute the loss $\mathcal{L}_{total}$ following Eq. (17);
13:           Update parameters $\theta$ with dataset $\mathcal{D}^{rs}$;
14:         **end if**
15:     **end for**
16: **end for**

| Model | KS | UCF-51 | Food-101 |
| --- | --- | --- | --- |
| Random re-sample | 60.78 | 68.59 | 84.21 |
| Inverse re-sample | 57.24 | 66.19 | 78.26 |
| DSR ($k = -0.5$) | 63.68 | 72.49 | 89.76 |
| DSR ($k = -1.0$) | 64.33 | 74.11 | 91.28 |
| DSR ($k = -1.5$) | 66.03 | 73.58 | 91.74 |
| DSR ($k = -2.0$) | 66.52 | 75.60 | 93.36 |
| DSR ($k = -2.5$) | 66.35 | 76.27 | 93.55 |
| DSR ($k = -3.0$) | 67.41 | 76.83 | 93.69 |

Table 5: Comparison with different re-sample frequencies. $k$ represents the slope of the sampling function $\mathcal{F}_s$, where red and blue indicates the best/runner-up performance.

a sample, $k$ is the slope of function and $m$ is the number of modality types, $round <>$ represents rounding operation.

Table 5 compares various re-sampling strategies, focusing on our proposed DSR method with different slopes $k$ for the sampling function $\mathcal{F}_s$. The slope $k$ determines the sampling frequency, where smaller $k$ values indicate higher re-sampling frequencies. In terms of performance, DSR with varying $k$ values consistently outperforms random and inverse re-sampling strategies. On the KS dataset, the best performance is achieved with $k = -3.0$, yielding an accuracy of 67.41%, which is 6.63% higher than the random re-sampling baseline. Similarly, for UCF-51 and Food-101, the optimal $k$ values lead to accuracy improvements of 8.24% and 9.48%, respectively. The results highlight the advantages of our DSR method, particularly its ability to dynamically adjust sampling based on modality contribution. DSR effectively balances the sampling frequency, ensuring that critical samples are revisited more frequently while avoiding over-sampling less informative samples. This dynamic

| DFF | BMML | DSR | KS | UCF-51 | Food-101 |
| --- | --- | --- | --- | --- | --- |
| | | | 59.61 | 68.23 | 82.38 |
| ✓ | | | 64.34 | 72.84 | 91.29 |
| | ✓ | | 63.65 | 72.68 | 90.83 |
| | | ✓ | 63.78 | 71.81 | 90.57 |
| ✓ | ✓ | | 65.09 | 74.29 | 92.40 |
| ✓ | | ✓ | 65.21 | 73.05 | 91.87 |
| | ✓ | ✓ | 64.58 | 73.76 | 92.16 |
| ✓ | ✓ | ✓ | 66.52 | 75.60 | 93.36 |

Table 6: Ablation study of each component in ARM.

approach allows for more efficient learning, leading to substantial performance improvements.

However, as the sampling frequency increases (i.e., as $k$ decreases), the computational cost also grows. This is because more frequent re-sampling leads to higher data processing requirements, which may limit the scalability of the approach in large-scale applications. Additionally, while increased sampling frequency can boost performance, there is a performance ceiling. For example, in UCF-51 and Food-101, the performance gains plateau as $k$ decreases from $-2.0$ to $-3.0$. This indicates that beyond a certain point, further increasing the sampling frequency yields diminishing returns. In summary, DSR provides a robust and flexible re-sampling strategy that outperforms traditional methods by dynamically adjusting to modality importance. However, it is important to balance the trade-off between sampling frequency and computational efficiency, as well as to recognize that the performance gains have practical limits.

**More Ablation Study**

We conducted ablation studies on the three modules in ARM, i.e. dynamic feature-level fusion (DFF), balanced min-max loss (BMML) and dynamic sample-level re-sample (DSR), to answer **Q6:** *How much does each module contribute in ARM?*

Table 6 presents the ablation study results, which highlight the contributions of each ARM component. The individual effects of DFF, BMML, and DSR demonstrate how each component impacts performance. When DFF is applied alone, the accuracy on the KS dataset improves from 59.61% to 64.34%, and similar gains are observed on UCF-51 and Food-101 datasets. DFF enhances feature interactions by dynamically adjusting the contributions of different modalities, allowing better feature fusion and synergy. This improvement illustrates how capturing the complementary information between modalities boosts overall performance. Incorporating BMML further boosts accuracy. All dataset sees an increase. BMML mitigates the impact of imbalanced contributions by balancing the influence of dominant and weaker modalities, which is crucial in real-world scenarios where some modalities may naturally dominate. The addition of DSR produces significant performance gains across all datasets. For example, on UCF-51, the accuracy jumps to 74.29%, while Food-101 reaches 92.40%. DSR dynamically adjusts the sampling frequency based on each modal-

| Model | KS | UCF-51 | Food-101 |
|---|---|---|---|
| Concatenation | 59.61 | 68.23 | 82.38 |
| Concatenation-ARM | 66.52 $\Delta$6.91 | 75.60 $\Delta$7.37 | 93.36 $\Delta$10.98 |
| Summation | 59.53 | 67.62 | 82.63 |
| Summation-ARM | 66.03 $\Delta$6.50 | 74.12 $\Delta$6.50 | 93.88 $\Delta$11.25 |
| MMTM | 63.92 | 70.21 | 90.63 |
| MMTM-ARM | 67.43 $\Delta$3.51 | 75.92 $\Delta$5.71 | 94.69 $\Delta$4.06 |
| CentralNet | 64.58 | 72.21 | 90.31 |
| CentralNet-ARM | 68.78 $\Delta$4.20 | 76.30 $\Delta$4.09 | 94.85 $\Delta$4.54 |

Table 7: Results of using ARM on various multimodal fusion Methods, $\Delta$ is accuracy enhancement.

ity's marginal contribution, ensuring that underrepresented modalities receive adequate focus during training. This dynamic resampling mechanism enhances model robustness, particularly when data distribution is skewed or modalities vary in importance.

When all three components: DFF, BMML, and DSR are integrated, the model achieves the highest performance across all datasets. The combined benefits stem from comprehensive enhancement strategies: better feature fusion, balanced contribution, and dynamic sample reweighting. The performance boost demonstrates that the integration of these components synergistically addresses the challenges of modality imbalance, feature misalignment, and suboptimal sampling. In summary, each ARM component provides distinct advantages, with DFF improving feature alignment, BMML addressing modality imbalance, and DSR optimizing sampling. Their combined impact leads to superior accuracy, reflecting their effectiveness in enhancing multimodal learning.

## Results on Other Fusion Methods

Notably, our method is not limited to fixed imbalanced multimodal learning frameworks; it can also be integrated into other existing approaches. In this section, we answer **Q7:** *how our model improves the performance of other multimodal fusion learning frameworks?*

Table 7 illustrates the performance improvement achieved by integrating ARM into various multimodal frameworks, including MMTM, and CentralNet. For the KS dataset, ARM consistently boosts accuracy across all frameworks, with improvements ranging from 6.91% for Concatenation to 4.20% for CentralNet. On other dataset, ARM shows similar trends. The substantial performance gains can be attributed to ARM's design, which dynamically balances the contributions of each modality. By addressing cross-modal biases, ARM prevents the dominance of any single modality and ensures more holistic learning. Additionally, ARM's focus on multimodal fusion allows it to capture complex relationships, effectively leveraging information across modalities. The consistent improvements across different architectures validate ARM's robustness and demonstrates its ability to adapt to various multimodal scenarios.

| Model | KS | UCF-51 | Food-101 |
|---|---|---|---|
| Local MI (MICCAI 2021) | 61.25 | 69.82 | 85.77 |
| MI-Dependency (EMNLP 2021) | 59.83 | 67.43 | 85.24 |
| Infomax (EMNLP 2021) | 62.54 | 70.04 | 89.26 |
| AMID (CVPR 2023) | 64.73 | 73.80 | 90.18 |
| ARM | 66.52 | 75.60 | 93.36 |

Table 8: Comparison with Mutual information (MI)-based multimodal learning methods.



Figure 6: The per-class accuracy (%) of recognition on UCF-51 comparing **ARM** with AMID.

Visualizations in Fig. 7 compare the feature space distributions of MMTM and CentralNet, before and after integrating our proposed ARM method. Without ARM, both MMTM and CentralNet display noticeable overlaps between clusters, indicating less discriminative feature spaces. After incorporating ARM, the cluster separations become more distinct and well-defined, suggesting improved feature representation and class separability. Specifically, ARM reduces intra-class variance and enhances inter-class separability, resulting in more cohesive clusters with minimal scatter. This enhancement translates into better classification performance. By addressing the imbalances in multimodal fusion, ARM not only strengthens the learning process for underrepresented modalities but also refines the overall joint feature space, leading to superior cluster organization. These qualitative improvements highlight how ARM effectively amplifies the strengths of existing multimodal networks like MMTM and CentralNet, demonstrating its general applicability and effectiveness in diverse multimodal scenarios.

## More Comparisons

In this section, we answer **Q8:** *How much does our method improve performance compared to existing MI-based multimodal learning methods?*

Table 8 presents a performance compared with other mutual information (MI)-based multimodal learning methods, including Local-MI (Liao et al. 2021), Mutual-Dependency (Colombo et al. 2021), Multimodal-Infomax (Han, Chen, and Poria 2021) and AMID (Chen et al. 2023). Our ARM outperforms the competing approaches by a significant margin on all three datasets. The notable performance gap highlights the effectiveness of ARM in capturing and balancing the complex multimodal relationships, which are critical

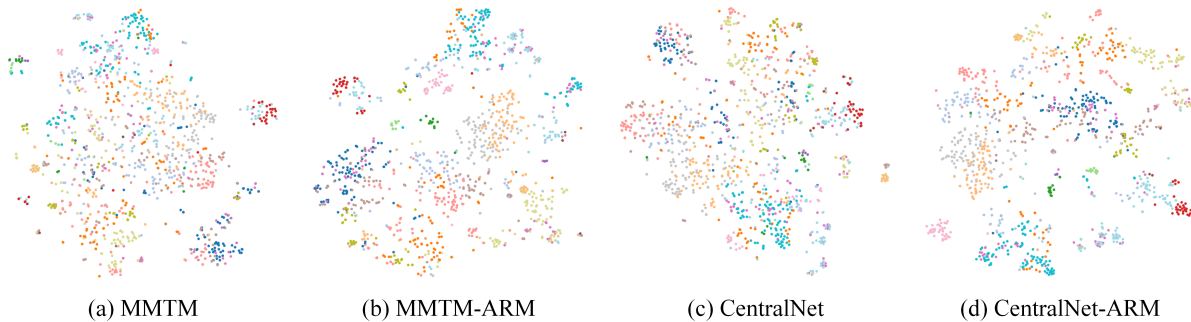(a) MMTM      (b) MMTM-ARM      (c) CentralNet      (d) CentralNet-ARM

Figure 7: Visual feature distribution of MMTM, MMTM-ARM and CentralNet, CentralNet-ARM visualized by t-SNE (Van der Maaten and Hinton 2008) on Kinetics Sounds dataset. Categories are indicated in different colors.
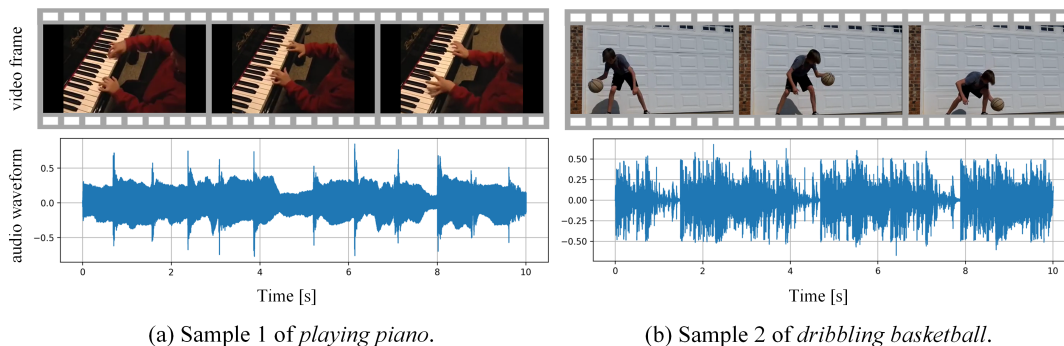


(a) Sample 1 of *playing piano*.      (b) Sample 2 of *dribbling basketball*.

Figure 8: Visualization of Audio-visual samples from Kinetics Sounds dataset.

for robust feature integration. The superior performance of ARM can be attributed to dynamically adjust contributions from each modality, allowing ARM address modality imbalance and information redundancy better, which are common issues in MI-based methods.

The per-class accuracy comparison between ARM and AMID in Fig. 6 demonstrates ARM's consistent performance advantage across various fine-grained action categories. While both methods show competitive results, ARM achieves superior accuracy in a majority of categories. This indicates that ARM is more effective in capturing the subtle nuances and intricate features within multimodal inputs, leading to better classification outcomes. Additionally, in categories where AMID struggles with lower accuracy, such as *Hand stand Pushups* and *HammerThrow*, ARM maintains a stable and high performance, reflecting its robustness in handling challenging and diverse actions. This consistency across the spectrum suggests that ARM effectively mitigates modality imbalance and enhances joint learning across different classes. Overall, the detailed analysis of fine-grained categories highlights ARM's strength in generalizing across varied scenarios while delivering more balanced and reliable results than AMID.

## Case Analysis of Modality Contribution

Here we provide a visualization instance to answer **Q9:** *How does ARM balance two modalities in samples with different contributions?*

Fig. 8 show two audio-visual multimodal pair of *playing piano* and *dribbling basketball* category, respectively. In Sample 1, the clear piano sound in the audio modality is easily recognizable, while the bouncing basketball action in Sample 2 is hard to detect due to unrelated background music interference. This could drag the joint contribution of all modalities by the more challenging-to-learn modality.

Fig. 9 compare contribution improvement for this two audio-visual samples under different imbalanced multimodal learning methods: Greedy, Sample-valuation, and our ARM. The contribution of each modality is tracked across training epochs (10, 40, 80). The results show that both Greedy and Sample-valuation exhibit fluctuating and imbalanced contributions, the focus on lower-contribution modality often leads to a decrease in higher-contribution modality. Meanwhile, they fail to maintain consistent contributions across epochs, resulting in instability. In contrast, our ARM demonstrates balanced and stable contribution enhancement from all modalities. The key reason for this advantage lies in ARM's dynamic re-sampling strategy and balanced min-max loss, which adaptively adjust the sampling frequency and contribution of each modality based on their marginal and joint contributions. This ensures that neither modality is overemphasized or ignored, leading to better generalization and more robust feature fusion. Consequently, ARM is able to achieve superior performance in scenarios with imbalanced modalities by maintaining consistent contribution levels, helping in maximizing the joint contribution gain.
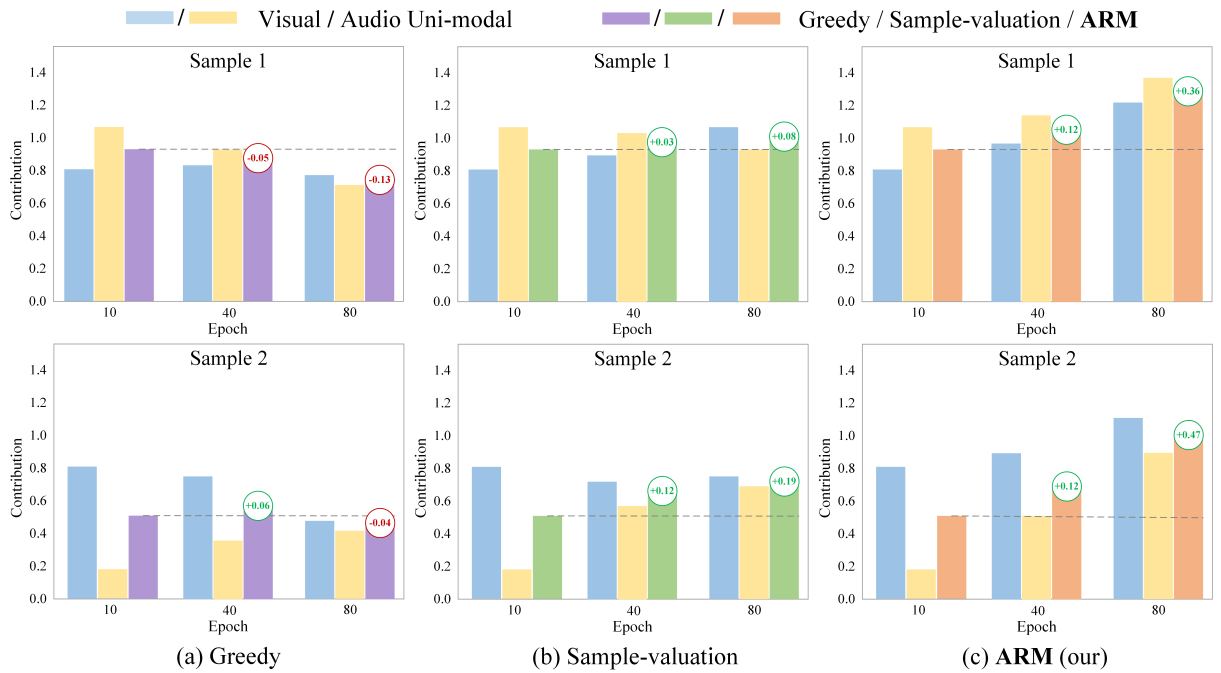
Figure 9: Contribution improvement compared. Other imbalanced multimodal learning methods: Greedy (Wu et al. 2022), Sample-valuation (Wei et al. 2024).