# MILES: Modality-Informed Learning Rate Scheduler for Balancing Multimodal Learning

Alejandro Guerra-Manzanares and Farah E. Shamout
Division of Engineering
New York University Abu Dhabi, Abu Dhabi, UAE
Email: {alejandro.guerra, farah.shamout}@nyu.edu

*Abstract*—The aim of multimodal neural networks is to combine diverse data sources, referred to as modalities, to achieve enhanced performance compared to relying on a single modality. However, training of multimodal networks is typically hindered by modality overfitting, where the network relies excessively on one of the available modalities. This often yields sub-optimal performance, hindering the potential of multimodal learning and resulting in marginal improvements relative to unimodal models. In this work, we present the Modality-Informed Learning ratE Scheduler (MILES) for training multimodal joint fusion models in a balanced manner. MILES leverages the differences in modality-wise conditional utilization rates during training to effectively balance multimodal learning. The learning rate is dynamically adjusted during training to balance the speed of learning from each modality by the multimodal model, aiming for enhanced performance in both multimodal and unimodal predictions. We extensively evaluate MILES on four multimodal joint fusion tasks and compare its performance to seven state-of-the-art baselines. Our results show that MILES outperforms all baselines across all tasks and fusion methods considered in our study, effectively balancing modality usage during training. This results in improved multimodal performance and stronger modality encoders, which can be leveraged when dealing with unimodal samples or absent modalities. Overall, our work highlights the impact of balancing multimodal learning on improving model performance.

*Index Terms*—multimodal learning, modality overfitting, learning rate scheduler, balanced training, joint fusion networks

## I. Introduction

Real-world prediction tasks are often multimodal in nature. From everyday activities like interpersonal communication [1], learning and cognition [2], to specialized tasks such as stock investment [3] and clinical differential diagnosis [4], humans integrate multiple sources of information to make accurate decisions and predictions. Despite this, most available datasets and machine learning models are unimodal, such that they rely on a single source of information, or modality, for a given prediction task. This dominant modeling approach overlooks the potential performance enhancements that could be achieved by combining multiple input modalities.

Multimodal machine learning seeks to enhance model performance by integrating related information from multiple sources. While several multimodal learning approaches have been successfully developed for specific problems [5], [6], multimodal training presents inherent challenges. These challenges include the need for effective data fusion, processing heterogeneous data modalities, and ensuring that the model can leverage the complementary nature of the multimodal data without being hindered by noise or inconsistencies. Furthermore, the alignment and synchronization of disparate data sources can be complex, especially when dealing with heterogeneous datasets that vary in format, scale, and quality.

Another significant challenge in multimodal machine learning is balancing the contributions of each modality to the learning process. Naive multimodal training leads to modality competition, where one modality becomes more dominant [7]. The multimodal model would then overfit to the dominant modality while underutilizing the remaining available modalities [7]. As an example, Fig. 1 shows naive runs for unimodal and multimodal models trained on the S-MNIST dataset. The left figure shows the training accuracy of the unimodal and multimodal encoders within a multimodal model, while the right figure shows validation accuracy of unimodal models trained independently. We can observe that while the image model has potential to reach a similar accuracy to that of the audio (right figure), it does not reach its full potential in the multimodal model (left figure), showing significantly lower performance. On the other hand, the audio encoder achieves comparable performance in both the unimodal and multimodal settings. The multimodal model (in blue) also tends to overfit to the audio model. Over-reliance on one modality may lead to sub-optimal performance, especially if that modality is noisy or less informative in certain contexts. Conversely, underutilization of a critical modality could result in missed opportunities to improve prediction accuracy. Additionally, the computational cost and model complexity tend to increase with the number of modalities considered, posing further obstacles to efficient model training and deployment.

To address the challenge of training multimodal networks in a balanced manner, we propose the Modality-Informed Learning ratE Scheduler (MILES). MILES is a learning rate scheduler for multimodal networks that leverages the conditional utilization rate per modality during training to dynamically adjust the learning rate, effectively balancing multimodal learning, increasing the contribution of the underutilized modalities and enhancing multimodal prediction performance. The main contributions of this work are summarized below:

- We present the Modality-Informed Learning RatE Scheduler, MILES, a novel learning rate scheduler for multimodal training that dynamically adjusts modality-specific learning rates based on each modality's contribution dur-
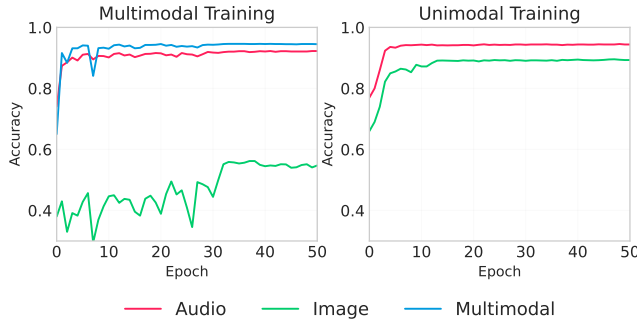
Fig. 1. Validation accuracy of multimodal model and its unimodal encoders (left) and unimodal models trained independently (right) using the S-MNIST dataset.

ing training, effectively balancing multimodal learning.

- We demonstrate that by increasing the utilization of non-dominant modalities and preventing overfitting to the dominant modality, MILES not only produces better multimodal results but also strengthens encoders performance, achieving accuracy comparable to unimodal models.
- We conduct extensive validation of MILES on four multimodal datasets using two fusion methods, demonstrating its superior performance compared to seven state-of-the-art approaches.

## II. RELATED WORK

### A. Multimodal joint fusion

Data fusion involves combining information from multiple sources (modalities) to extract complementary insights, leading to more comprehensive and better-performing machine learning models than those relying on a single data source. The three primary fusion strategies are early fusion, joint fusion, and late fusion [8]. Early fusion, also known as feature-level fusion, combines the available modalities into a single feature vector, which is used as input for a machine learning model. On the other hand, late fusion, also known as decision-level fusion, aggregates predictions from multiple models to compute the final prediction. Finally, joint fusion, also known as intermediate fusion, involves combining learned feature representations from intermediate layers of neural networks (encoders), for further processing. Joint fusion models are trained end-to-end, allowing feature representations to improve with each iteration. Our work focuses on multimodal learning in joint fusion networks, which is the most common setting of multimodal fusion. These networks combine features at different levels of the model, leading to a more comprehensive understanding of the data and generally better overall performance. We demonstrate the impact of our proposed approach for the most widely used fusion methods: feature concatenation and feature summation [8].

### B. Modality competition

Even though multimodal fusion models have the potential to outperform unimodal models, the joint training of multimodal networks often fails [9] or produces a sub-optimal model that does not surpass the best unimodal model [10]. During joint training, available modalities compete with each other, leading to a situation where only a small subset of modalities are effectively learned while the others are neglected and remain underexplored [9]. In practice, as shown in previous work, only one of the available modalities dominates the learning process [7]. Possible reasons behind the tendency of multimodal models to overfit to one modality while neglecting others are better alignment of certain modalities with the encoding network's random initialization [9], faster learning speeds of specific modalities [7], and differences in convergence rates and generalization performance among modalities [10], [11].

### C. Balanced multimodal learning

To address modality competition, avoid overfitting, and balance modality contributions, several methods have been proposed for end-to-end training of multimodal fusion networks. [11] introduced the Modality-Specific Early Stopping (MSES) method, which tackles the differing convergence rates and generalization performance of various modalities by applying a global learning rate to the network and using early stopping during joint training for modalities that have converged based on validation performance. The training continues until all parts of the network have converged.

Building on similar reasoning, [12] proposed the Modality-Specific Learning Rate (MSLR) method. MSLR is based on the observation that while a global learning rate may work well for some modalities (and lead to overfitting), it may be outside the effective range for others (causing them to be ignored). To address this, MSLR advocates using distinct learning rates for the multimodal model and each encoder. The proposed MSLR variants include: MSLR-K (the *keep* strategy, which maintains the best fine-tuned unimodal learning rates for different modalities fixed during multimodal model training), MSLR-S (the *smooth* strategy, which adjusts the learning rates of different modalities to be closer to the average learning rate across modalities), and MSLR-D (the *dynamic* strategy, which scales the learning rate for each unimodal encoder during training by a small factor based on the unimodal average prediction performance on the validation set over the last several epochs).

[10] proposed Gradient Blending (G-Blend), which directly modifies the gradient descent process by altering the loss function to a weighted sum of multiple unimodal losses. The weight for each modality is computed based on an Overfitting-to-Generalization Ratio (OGR), which describes the overfitting conditions for each modality. To compute OGR, each unimodal model is trained individually for the first several epochs, using the same learning rate for both the modality-specific and multimodal models. [13] introduced the On-the-fly Gradient Modulation (OGM) method to adaptively control the optimization of each modality, via monitoring the discrepancy of their contribution towards the learning objective. OGM-GE refines this method by incorporating dynamic Gaussian noise

during training to prevent a generalization drop from gradient modulation.

Compared to existing approaches, MILES computes the conditional utilization rate per modality after each epoch as a proxy for each modality's contribution to the multimodal model and adjusts the learning rate only for the dominant modality accordingly. This simple yet effective adjustment does not require calculating any weights, performing uni-modal pretraining, or setting individual learning rates, which necessitate increased computational resources and hyper-parameter tuning rounds for offline unimodal model pretraining prior multimodal training. MILES only requires a global learning rate, as the modified learning rate is derived from the global learning rate using a scaling factor. All of the aforementioned state-of-the-art approaches are used as baselines in our work.

## III. METHODOLOGY

In this section, we begin by introducing some preliminaries and the conditional utilization rate, a key component leveraged by our approach to dynamically adjust the learning rate during training. Following that, we present our proposed method, the Modality-Informed Learning ratE Scheduler (MILES), designed to address modality competition and balance multimodal utilization during the training process.

### A. Preliminaries

In this work, we assume the presence of two modalities, $A$ and $B$, with the goal of predicting a target variable $\mathbf{y}$. Each modality is processed by its own encoder within a multimodal network. Let $\mathbf{A} \in \mathbb{R}^{n \times d_A}$ and $\mathbf{B} \in \mathbb{R}^{n \times d_B}$ represent the feature matrices for modalities $A$ and $B$, respectively, where $n$ is the number of samples, $d_A$ is the dimensionality of modality $A$, and $d_B$ is the dimensionality of modality $B$. The multimodal network consists of two modality encoders, $f_A$ and $f_B$, which transform the input features into encoded representations:

$$\mathbf{z}_A = f_A(\mathbf{A}), \quad \mathbf{z}_B = f_B(\mathbf{B}), \quad (1)$$

where $\mathbf{z}_A$ and $\mathbf{z}_B$ are the latent representations of modalities $A$ and $B$, respectively. The goal is to predict the target variable $\mathbf{y}$ using a function $g$ that combines the encoded features:

$$\hat{\mathbf{y}}_{AB} = g_{AB}(\mathbf{z}_A, \mathbf{z}_B), \quad (2)$$

where $\hat{\mathbf{y}}_{AB}$ is the prediction using the multimodal network. The prediction is learned by optimizing the network parameters of $f_A$, $f_B$ and $g_{AB}$ jointly to minimize the loss function:

$$\mathcal{L} = \mathcal{L}_{AB}(\hat{\mathbf{y}}_{AB}, \mathbf{y}) + \mathcal{L}_A(\hat{\mathbf{y}}_A, \mathbf{y}) + \mathcal{L}_B(\hat{\mathbf{y}}_B, \mathbf{y}), \quad (3)$$

where $\mathcal{L}_{AB}(\hat{\mathbf{y}}_{AB}, \mathbf{y})$, is the loss for the multimodal model, $\mathcal{L}_A(\hat{\mathbf{y}}_A, \mathbf{y})$ refers to the loss of the encoder of modality A, $\mathcal{L}_B(\hat{\mathbf{y}}_B, \mathbf{y})$ for the modality B encoder. $\hat{\mathbf{y}}_i$ refers to the predicted label by the corresponding model and $\mathbf{y}$ is the true label.

### B. Motivation

The *greedy learner hypothesis*, introduced by [7], establishes that a multimodal learning process where the network is trained to minimize modality-specific losses is inherently greedy. This process tends to rely on only one of the available input modalities — the one that is the fastest to learn from. However, as shown in Figure 1 and in related work such as by [10], this phenomenon is also observed in joint fusion networks, which are typically trained to minimize a single multimodal loss (the $\mathcal{L}_{AB}(\hat{\mathbf{y}}_{AB}, \mathbf{y})$ term in Equation 3) or a sum of multimodal and modality-specific losses (Equation 3). To address this problem in joint fusion networks, we propose MILES, a method designed to regulate the learning speed at which modalities are learned by the multimodal model. MILES dynamically adjusts the modality-specific learning rates based on an epoch-wise estimation of each modality's marginal contribution to the overall multimodal performance. This adjustment slows down the learning of the dominant modality, allowing the non-dominant modalities to be learned more effectively.

### C. Conditional utilization rate

We redefine and repurpose the *conditional utilization rate*, originally presented by [7] and modified by [14], denoted as $\mathbf{u}$, to characterize epoch-wise modality usage in multimodal joint fusion networks. We define $\mathbf{u}$ for multimodal joint fusion networks as follows:

$$\mathbf{u_A} = \frac{M(\hat{\mathbf{y}}_{AB}) - M(\hat{\mathbf{y}}_B)}{M(\hat{\mathbf{y}}_{AB})}, \mathbf{u_B} = \frac{M(\hat{\mathbf{y}}_{AB}) - M(\hat{\mathbf{y}}_A)}{M(\hat{\mathbf{y}}_{AB})}, \quad (4)$$

where $M(\cdot)$ represents the performance metric used (e.g., accuracy, F1 score), $\mathbf{u_A}$ determines the conditional utilization rate for modality $A$, and $\mathbf{u_B}$ for modality $B$. Intuitively, it quantifies the marginal contribution of a specific modality to the fusion model's performance.

We define $\delta_{AB}$ as the absolute value of the difference between conditional utilization rates:

$$\delta_{AB} = |\mathbf{u_A} - \mathbf{u_B}|. \quad (5)$$

It enables the evaluation of imbalanced modality usage within the multimodal fusion model. Note that $\delta_{AB} \in \mathbb{R}$ : $s.t.\ 0 \leq \delta_{AB} \leq 1$, with values closer to one indicating imbalanced modality usage and modality overfit, since the difference in conditional utilization rates would be high indicating that one modality is being used more than the other.

### D. Modality-Informed Learning ratE Scheduler (MILES)

Our proposed methodology for addressing modality competition and enhancing multimodal training, is described in Algorithm 1. MILES only requires a global learning rate to be specified, and it re-adjusts the learning rates of individual modalities by scaling them with a factor of $\mu$ when certain conditions are met. These conditions are governed by the target difference threshold, denoted as $\tau$. By combining these two hyper-parameters ($\tau$ and $\mu$) following the conditional

**Algorithm 1** MILES

**Input**: Target difference threshold: $\tau$, Training epochs: $N$, Reduction factor: $\mu$, Global learning rate: $\alpha$

1: Initialize variables: $\delta_{AB} \leftarrow 0$, $\alpha_{AB} \leftarrow \alpha$, $\alpha_A \leftarrow \alpha$, $\alpha_B \leftarrow \alpha$
2: **for** $i = 1, \ldots, N$ **do**
3:     Train_epoch($i$)
4:     Validation_epoch($i$)

5:     Compute $\mathbf{u_A}$ {$\mathbf{u}$ of modality A as in Eq. 4}
6:     Compute $\mathbf{u_B}$ {$\mathbf{u}$ of modality B as in Eq. 4}
7:     Compute $\delta_{AB}$ {Difference between modalities as in Eq. 5}

8:     **if** $(\delta_{AB} \leq \tau) \vee (\mathbf{u_A} < 0 \wedge \mathbf{u_B} < 0)$ **then**
9:         $\alpha_A \leftarrow \alpha_{AB}$
10:        $\alpha_B \leftarrow \alpha_{AB}$
11:    **else**
12:        **if** $\mathbf{u_A} > 0 \wedge \mathbf{u_B} < 0$ **then**
13:            $\alpha_A \leftarrow \mu \cdot \alpha_{AB}$
14:        **else if** $\mathbf{u_A} < 0 \wedge \mathbf{u_B} > 0$ **then**
15:            $\alpha_B \leftarrow \mu \cdot \alpha_{AB}$
16:        **else**
17:            **if** $\mathbf{u_A} < \mathbf{u_B}$ **then**
18:                $\alpha_B \leftarrow \mu \cdot \alpha_{AB}$
19:            **else**
20:                $\alpha_A \leftarrow \mu \cdot \alpha_{AB}$
21:            **end if**
22:        **end if**
23:    **end if**
24: **end for**

statements and computations defined in Algorithm 1, MILES dynamically adjusts the learning rate per modality ($\alpha_A$ and $\alpha_B$) during training based on differences in conditional utilization rates.

After every training epoch and subsequent validation stage, MILES is applied to tune the learning rates per modality for the next training iteration (lines 3-4, Algorithm 1). The first step of the MILES procedure involves calculating the conditional utilizations per modality ($\mathbf{u_{A,B}}$) using Eq. 4, and the difference between them ($\delta_{AB}$) using Eq. 5 (lines 5-7, Algorithm 1). Based on computed values $\mathbf{u_{A,B}}$, $\delta_{A,B}$ and set hyper-parameters $\tau$ (target difference threshold) and $\mu$ (reduction factor), the following conditions apply (lines 8-24):

- If $\delta_{AB}$ is below or equal to the target threshold $\tau$, or $\mathbf{u_A} < 0$ and $\mathbf{u_B} < 0$, MILES takes no action. In the former case, the target threshold is met so no action from MILES is required. In the latter case, the unimodal encoders outperform the multimodal model, which typically occurs during the first few epochs of training when the multimodal model has not yet started learning effectively. In such instances, it is preferable to wait until the multimodal model begins to learn, so MILES takes no action. In either case, the learning rates for both modalities rate are set as the global learning rate for the next training epoch (lines 8-10, Algorithm 1).
- If the previous condition is not satisfied, meaning $\delta_{AB} > \tau$, the following conditions and actions apply:
  - If $\mathbf{u_A} > 0$, $\mathbf{u_B} < 0$, and $\delta_{AB} > 0$ then this implies that modality B is under-utilized and the learning

rate of modality A is scaled by $\mu$ for the next epoch (lines 12-13). The aim is to slow down the learning of modality A, the over-utilized modality, allowing the model to learn from modality B more effectively.
  - If $\mathbf{u_A} < 0$, $\mathbf{u_B} > 0$, and $\delta_{AB} > 0$ then this implies that modality A is under-utilized and the learning rate of modality B is scaled by $\mu$ for the next epoch (lines 14-15). This is the inverse of the previous condition, where the same rationale applies.
  - If $\mathbf{u_A} > 0$ and $\mathbf{u_B} > 0$, then this indicates that the multimodal model is outperforming both modality encoders, suggesting that the multimodal model is learning effectively. However, to assess if there is modality overfitting, $\mathbf{u_A}$ and $\mathbf{u_B}$ need to be further compared:
    * If $\mathbf{u_A} < \mathbf{u_B}$ then modality A is under-utilized and the learning rate of modality B is scaled by $\mu$ for the next epoch (lines 17-18).
    * If the opposite is true, $\mathbf{u_A} > \mathbf{u_B}$, then modality B is under-utilized and the learning rate of modality A is scaled by $\mu$ for the next epoch (lines 19-20).

Note that the following assumptions also apply:
- $\tau \in [0.0, 1.0]$ because $\tau$ only takes possible values of $\delta_{AB}$, and the range of $\delta_{AB}$ (Eq. 5) lies within the interval $[0.0, 1.0]$.
- $\mu \in \left\{ \frac{a}{b} \,\middle|\, a = 1, b \in \mathbb{R}, b \neq 0 \right\}$ as the usage of $\mu$ in Algorithm 1 involves multiplying it by the current learning rate to effectively reduce it, $\mu$ cannot be greater than 1.
- An auxiliary supervised loss term is added for each encoder to the objective function (e.g., for modality A, a binary cross entropy loss term, $\mathcal{L}_A(\mathbf{y}, \hat{\mathbf{y}}_A)$, would be added, where $\mathbf{y}$ refers to the ground-truth labels and $\hat{\mathbf{y}}_A$ are the unimodal encoder predictions).
- Only one modality's learning rate is adjusted per epoch.

## IV. EXPERIMENTS

In this section, we provide an overview of the datasets, baselines and experimental settings used in the main experiments and ablation studies. We utilize various multimodal datasets, used as benchmarks in prior work, either in their unimodal or multimodal versions. For reproducibility, we make our implementation publicly available at https://github.com/nyuad-cai/MILES.

### A. Datasets

**CREMA-D** [15] is an audio-visual dataset for multimodal emotion recognition. It includes 7,442 video clips from 91 actors speaking a selection of 12 sentences. The utterances express one of six common emotions: anger, happiness, disgust, fear, neutral, and sadness. We split the full dataset into training, validation, and test sets with a 70-15-15 split. The training set includes 5,210 samples, while both the validation and test sets include 1,116 samples.

**S-MNIST** [16] is an audio-visual dataset designed for benchmarking multimodal classification. It pairs the original MNIST [17] with a spoken digits database from Google

| DATASET | Modality A | Performance | Modality B | Performance |
|---------|-----------|-------------|-----------|-------------|
| S-MNIST | Audio | 88.9 A. | Image | 99.0 A. |
| CREMA-D | Audio | 61.9 A. | Video | 60.7 A. |
| MM-IMDb | Text | 64.4 A. | Image | 38.9 A. |
| LUMA | Audio | 97.2 F1 | Image | 79.1 F1 |

Speech Commands [18]. We sampled 8,000 and 2,000 samples from the original training set as our training and validation sets. We sampled 2,000 samples from the original test set for our test set.

**LUMA** [19] includes a multimodal image-audio dataset for benchmarking multimodal learning. It contains images from a 50-class subset of CIFAR-10/100 and class label utterances. The original dataset is imbalanced, with the top 22 classes having 1,500 or more samples each. We generate a balanced set with 1,500 instances per class, totaling 33,000 pairs, which are split into 25,000 training, 4,000 validation, and 4,000 test samples.

**MM-IMDb** [20] is the largest publicly available multimodal dataset for movie genre multilabel classification. It contains plot summaries (text modality) and posters (image modality) for 25,959 movies. The dataset is imbalanced, with *drama* being the most prevalent category (13,967 samples) and *film-noir* the least prevalent (338 samples). The dataset is split into three subsets: the training, validation, and test sets contain 15,552, 2,608, and 7,799 samples, respectively.

### B. Model training and evaluation

**Architecture.** The *vanilla* multimodal neural network architecture employed in our experiments is composed of two unimodal networks, one for each modality. We parameterize the encoders for each modality as follows: (i) BERT [21] for the text modality, and (ii) ResNet [22] variants for the image and audio modalities. The intermediate representations from the encoders are then fused and used as input for a linear layer that outputs the multimodal prediction. We evaluate all tasks using two fusion methods: concatenation and summation.

**Baselines.** We compare MILES to seven state-of-the-art approaches for balanced multimodal learning:

- MSES [11] uses early stopping for balancing multimodal learning based on convergence and generalization performance.
- MSLR [12] uses the best unimodal learning rate to enhance modality learning. We evaluate all variants: (i) MSLR-K, (ii) MSLR-S, (iii) MSLR-D.
- OGM [13] adaptively controls the optimization of each modality and OGM-GE adds noise dynamically to OGM to avoid possible generalization drops.
- G-Blend [10] adds multiple unimodal losses to the loss function (similar to MILES, as shown in Equation 3), and weights them based on the OGR of each modality.

**Implementation details**. We apply all of the baselines and our proposed approach on the vanilla architecture using the

two fusion variants. We use open source implementations if available. For the CREMA-D dataset, we extract one frame from each clip and process the audio data as a spectrogram of size $257 \times 299$ with window length of 512 and overlap of 353. We resize both modalities to $224 \times 224$ images (3 channels for image and one for audio) and use ResNet-18 [22] as backbones. We train all CREMA-D models for 200 epochs. For the S-MNIST dataset, we use the processed data provided by [16]. We reshape and resize both modalities to be $28 \times 28$ images (one channel) and use ResNet-10 as backbones. We train all S-MNIST models for 50 epochs. For the MM-IMDb dataset, we pre-process the data following [20]. We use a base BERT model with pre-trained weights, and a ResNet-50 model as text and image backbones, respectively. We train all these models for 50 epochs. For the LUMA dataset, we pre-process the data as provided by [19], generating a mel spectrogram (96 filterbanks) for the audio modality (shape $96 \times 96 \times 1$) and resizing CIFAR-10/100 to $96 \times 96 \times 3$. We use ResNet-6 as backbones and train all models for 60 epochs. For all experiments, we use the Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ . We train all models using an Nvidia A100 GPU.

**Hyper-parameter tuning.** For fair comparison of methods across models, we use random hyper-parameter tuning for all experiments. We randomly sample 80 combinations of hyper-parameters per model for each dataset, fusion method, and baseline using the recommended hyper-parameter ranges suggested by authors in the original work, or based on ranges determined according to our experiments. We select the best model for each method based on the multimodal performance on the validation set and evaluate each selected model on the final test set. We report the best performance results on the test set.

**Performance metrics.** For the S-MNIST, LUMA, and CREMA-D tasks, we report the best accuracy results on the test set, while for the highly imbalanced MM-IMDb, we report the best F1 score results on the test set. In addition, we report the difference between the unimodal classification heads for each method evaluated. Specifically, we subtract the weaker modality from the stronger modality on the best model on the test set. Formally, assuming $A$ is the stronger modality and $B$ is the weaker modality, we define it as $\Delta_{AB} = A_A(\cdot) - A_B(\cdot)$, where $A(\cdot)$ denotes the classification accuracy metric.

**Ablation studies.** To understand the impact of each hyper-parameter in the MILES algorithm and to provide recommendations for their use, we conduct three ablation studies:

1) $\tau$ **sensitivity analysis:** We vary the value of $\tau$ and analyze its impact on test set accuracy while keeping the other hyper-parameters fixed.
2) $\mu$ **sensitivity analysis:** We vary the value of $\mu$ and analyze its impact on test set accuracy while keeping the other hyper-parameters fixed.
3) **Computing u based on training metrics:** We apply MILES by computing the conditional utilization rates (Eq. 4) using the training set performance metrics instead of the validation set metrics.

TABLE II

ACCURACY RESULTS ON THE TEST SET FOR THE MULTIMODAL FUSION CONCATENATION HEAD (⊕, GREEN COLUMN) AND EACH MODALITY ENCODER. THE YELLOW COLUMN REPORTS RESULTS FOR THE OVERFITTING MODALITY WHILE THE BLUE COLUMN REPORTS THE RESULTS FOR THE UNDER-UTILIZED MODALITY. THE DIFFERENCES IN PERFORMANCE BETWEEN THE UNIMODAL CLASSIFICATION HEADS ($\Delta_{AB}$) ARE REPORTED IN THE WHITE COLUMN. THE BEST RESULTS ARE SHOWN IN BOLD.

| MODEL | CREMA-D | | | | S-MNIST | | | | LUMA | | | | MM-IMDb | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mathbf{A}_{A\oplus V}$ | $\mathbf{A}_A$ | $\mathbf{A}_V$ | $\Delta_{AV}$ | $\mathbf{A}_{I\oplus A}$ | $\mathbf{A}_I$ | $\mathbf{A}_A$ | $\Delta_{IA}$ | $\mathbf{A}_{A\oplus I}$ | $\mathbf{A}_A$ | $\mathbf{A}_I$ | $\Delta_{AI}$ | $\mathbf{F1}_{T\oplus I}$ | $\mathbf{F1}_T$ | $\mathbf{F1}_I$ | $\Delta_{TI}$ |
| Vanilla | 62.6 | 58.2 | 26.6 | +31.6 | 98.4 | 98.5 | 53.9 | +44.6 | 98.2 | 71.0 | 29.8 | +41.2 | 63.3 | 62.5 | 27.9 | +34.6 |
| + MSLR-K [12] | 65.1 | **59.9** | 29.7 | +30.2 | 99.0 | **98.8** | 58.9 | +39.9 | 98.7 | 87.8 | 38.1 | +49.7 | 63.4 | 62.4 | 27.9 | +34.5 |
| + MSLR-S [12] | 64.8 | 58.2 | 29.2 | +29.0 | 99.1 | 98.6 | 66.1 | +32.5 | 98.5 | 93.2 | 35.7 | +57.5 | 63.6 | 62.7 | 28.2 | +34.5 |
| + MSLR-D [12] | 64.3 | **59.9** | 39.7 | +20.2 | 99.0 | **98.8** | 76.9 | +21.9 | 98.7 | 93.3 | 37.0 | +56.3 | 63.6 | 62.6 | 28.5 | +34.1 |
| + OGM [13] | 63.9 | 57.9 | 34.6 | +23.3 | 98.9 | **98.8** | 64.5 | +34.3 | 98.5 | 86.4 | 44.0 | +42.4 | 63.7 | 63.0 | 28.3 | +34.7 |
| + OGM-GE [13] | 70.4 | 54.6 | 44.2 | +10.4 | 99.0 | 97.3 | 73.0 | +24.3 | 96.2 | 74.7 | 49.7 | +25.0 | 64.0 | 62.6 | 29.3 | +33.3 |
| + MSES [11] | 71.6 | 57.7 | 55.9 | +1.8 | 98.4 | 98.3 | 77.8 | +20.5 | 98.1 | 88.3 | 66.2 | +22.1 | 64.4 | 63.4 | 30.6 | +32.8 |
| + G-Blend [10] | 71.8 | 58.4 | 56.2 | +2.2 | 98.8 | 98.1 | 75.9 | +22.2 | 98.5 | 88.4 | 67.7 | +20.7 | 64.3 | 63.4 | 30.3 | +33.1 |
| + MILES (Ours) | **75.1** | **59.9** | **60.8** | **-0.9** | **99.8** | **98.8** | **84.9** | **+13.9** | **99.7** | **95.1** | **75.1** | **+20.4** | **65.1** | **64.2** | **36.6** | **+27.6** |

TABLE III

ACCURACY AND F1 PERFORMANCE RESULTS ON THE TEST SET FOR THE MULTIMODAL FUSION SUMMATION HEAD (+, GREEN COLUMN) AND EACH MODALITY ENCODER. THE YELLOW COLUMN REPORTS RESULTS FOR THE DOMINANT MODALITY WHILE THE BLUE COLUMN REPORTS THE RESULTS FOR THE UNDER-UTILIZED MODALITY. THE DIFFERENCES IN PERFORMANCE BETWEEN THE UNIMODAL CLASSIFICATION HEADS ($\Delta_{AB}$) ARE REPORTED IN THE WHITE COLUMN. THE BEST RESULTS ARE SHOWN IN BOLD.

| MODEL | CREMA-D | | | | S-MNIST | | | | LUMA | | | | MM-IMDb | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mathbf{A}_{A+V}$ | $\mathbf{A}_A$ | $\mathbf{A}_V$ | $\Delta_{AV}$ | $\mathbf{A}_{I+A}$ | $\mathbf{Acc}_I$ | $\mathbf{A}_A$ | $\Delta_{IA}$ | $\mathbf{A}_{A+I}$ | $\mathbf{A}_A$ | $\mathbf{A}_I$ | $\Delta_{AI}$ | $\mathbf{F1}_{T+I}$ | $\mathbf{F1}_T$ | $\mathbf{F1}_I$ | $\Delta_{TI}$ |
| Vanilla | 62.9 | 58.5 | 27.2 | +31.3 | 98.3 | 98.0 | 55.8 | +42.2 | 98.5 | 86.7 | 43.2 | +43.5 | 63.1 | 63.1 | 27.8 | +35.3 |
| + MSLR-K [12] | 65.7 | 57.4 | 33.3 | +24.1 | 98.7 | 98.5 | 68.2 | +30.3 | 98.8 | 87.0 | 43.6 | +43.4 | 63.5 | 61.5 | 29.1 | +32.4 |
| + MSLR-S [12] | 66.0 | 59.9 | 31.1 | +28.8 | 98.8 | **98.8** | 58.1 | +40.7 | 98.7 | 92.0 | 44.3 | +47.7 | 63.9 | 61.8 | 28.9 | +32.9 |
| + MSLR-D [12] | 65.4 | 59.9 | 34.9 | +25.0 | 99.1 | 98.7 | 73.5 | +25.2 | 98.9 | 91.8 | 48.2 | +43.6 | 64.0 | 62.0 | 29.3 | +32.7 |
| + OGM [13] | 65.6 | 56.8 | 36.9 | +19.9 | 98.6 | 97.8 | 61.8 | +36.0 | 98.7 | 87.3 | 46.7 | +40.6 | 63.9 | 62.2 | 28.7 | +33.5 |
| + OGM-GE [13] | 68.6 | 52.5 | 43.6 | +8.9 | 99.0 | 96.9 | 79.1 | +17.8 | 97.7 | 73.2 | 48.4 | +24.8 | 63.9 | 62.0 | 29.0 | +33.0 |
| + MSES [11] | 72.1 | 55.1 | 58.1 | -3.0 | 98.9 | 98.6 | 81.4 | +17.2 | 98.4 | 90.2 | 70.6 | +19.6 | 64.0 | 63.8 | 29.1 | +34.7 |
| + G-Blend [10] | 72.3 | 54.9 | 57.6 | -2.7 | 99.2 | **98.8** | 80.1 | +18.7 | 98.7 | 90.3 | 71.1 | +19.2 | 64.3 | 63.1 | 29.1 | +34.0 |
| + MILES (Ours) | **75.7** | **60.0** | **60.8** | **-0.8** | **99.8** | **98.8** | **85.0** | **+13.8** | **99.8** | **95.7** | **76.6** | **+19.1** | **65.0** | **64.0** | **33.1** | **+30.9** |

For the sensitivity analysis studies, we fix the learning rate to the one corresponding to the best results. For the third ablation study, we perform hyper-parameter tuning as described in the previous paragraph.

## V. RESULTS

This section provides the experimental results for unimodal baselines, multimodal performance and all ablation studies. Based on these, we derive some recommendations on how to use MILES more efficiently.

### A. Unimodal performance results

We report the results of the best unimodal models for each test set in Table I for comparison with multimodal model performance. As observed, the four datasets exhibit varying performance across modalities. Specifically, for CREMA-D, despite being a more complex task (emotion recognition) than the other three tasks, both modalities have similar unimodal accuracy (61.9% vs. 60.7%). The S-MNIST dataset is relatively easier for unimodal models, with image data achieving 99% accuracy and 88.9% accuracy for the audio modality. For the LUMA dataset, the audio modality outperforms the visual modality by a wide margin (97.2% vs. 79.1%). Finally, for the MM-IMDb task, multilabel movie genre classification, the text modality is significantly better than the image modality (64.4% vs. 38.9%).

### B. Multimodal performance results

**Concatenation**. Table II shows the results for all baseline models and MILES on the four datasets using feature concatenation as the fusion method. MILES consistently outperforms all baselines on multimodal and unimodal performance across datasets, with the greatest margin for the CREMA-D dataset, which has the lower baseline performance (62.6% vs. 75.1% for multimodal performance). Apart from enhancing multimodal prediction, MILES produces strong modality encoders with accuracy close to the unimodal baselines shown in Table I. The modality difference measure, $\Delta_{AB}$, reaches its minimum value for MILES in all tasks, indicating that MILES can effectively balance modality usage for both multimodal and unimodal predictions. More interestingly, for the CREMA-D task, the difference is negative, indicating that the weaker modality became the stronger for the MILES model, with the stronger modality still on par with the best baselines methods. Comparing the unimodal models (see Table I) and unimodal encoders (see Table II), on average, there is a gap of $\approx 1.8\pm1.6\%$ in performance between the best unimodal model and the top-performing MILES modality encoders (across the eight encoders), with the smallest gap being $0\%$ for CREMA-D video and the largest gap being $4\%$ for the LUMA image encoder. In all cases, MILES boosts the predictive performance of the non-dominant modality, especially for LUMA (29.8% vs. 75.1%), while keeping top accuracy metrics for the dominant modality.

**Summation**. Table III shows the results for all baseline models and MILES on the four datasets using feature summation as the fusion method. Similar to the setting using feature concatenation, MILES outperforms all baselines on multimodal and unimodal performance across all datasets.

While we observe the same overall trends as for concatenation, the results for the feature summation technique are greater than for concatenation for most models. The same observations as in the concatenation setting apply for the difference between encoders, $\Delta_{AB}$, with MILES providing the minimum values, thereby producing not only stronger but also more balanced encoders. Comparing the performance of unimodal models with the unimodal encoders of the multimodal networks, MILES consistently produces the best multimodal results across all datasets and also generates the strongest unimodal encoders, with an average performance gap between the best models and MILES encoders of $\approx 2.0 \pm 2.0\%$ (ranging from 0% for the CREMA-D video encoder to 5.8% for the MM-IMDb image encoder).

Overall, these results show that regardless of the fusion method employed and across datasets and diverse backbone architectures, MILES allows overcoming the modality overfit of conventional training, surpassing all state-of-the-art methods in enhancing both multimodal and unimodal performance, especially for the non-dominant modality.

### C. Ablation studies results

In the following paragraphs, we provide the results for the ablation studies, which were performed using the LUMA dataset. For both hyper-parameter sensitivity analyses, we use the best configurations for MILES in our experiments ($\tau = 0.2$ and $\mu = 0.5$).

$\tau$ **sensitivity analysis**. Table IV shows the results of varying $\tau$ while keeping the other hyper-parameters fixed (using the hyper-parameters of the best MILES model on the LUMA dataset from Table II and Table III). The results show that smaller $\tau$ values (e.g., $\tau = 0.0$) force the model to balance both modalities, boosting the performance of the non-dominant modality. However, this comes at the cost of notably lowering the performance of the dominant modality, which negatively impacts both unimodal and multimodal performance, resulting in improved but sub-optimal models.

The same occurs at the other end of the spectrum when $\tau$ is larger (e.g., $\tau = 0.5$), as it does not restrict the learning of the dominant modality, allowing the model to overfit. Nevertheless, the multimodal model improves the non-dominant modality performance via the added modality-specific losses (Equation 3), but the results are suboptimal with respect to the best MILES results.

$\mu$ **sensitivity analysis**. Table V reports the results of varying $\mu$ while keeping the other parameters fixed using the hyper-parameters of the best MILES model on the LUMA dataset (Table II and III). The results show that smaller values of $\mu$, such as 0.05 and 0.01, which imply a greater reduction in the learning rate of the dominant modality, significantly affect the learning of that specific modality. This impacts overall model performance, enhancing the non-dominant modality's performance but potentially degrading the dominant modality's performance. On the contrary, larger values closer to one may not enhance the learning of the non-dominant modality. The

| $\tau$ | Fusion | $\text{Acc}_{IA}$ | $\text{Acc}_A$ | $\text{Acc}_I$ |
|---|---|---|---|---|
| 0.5 | | 98.9 | 97.4 | 57.6 |
| 0.3 | | 99.1 | 95.4 | 66.7 |
| 0.2 | $\oplus$ | **99.7** | **95.1** | **75.1** |
| 0.1 | | 99.1 | 93.1 | 76.8 |
| 0.0 | | 98.7 | 91.7 | 77.5 |
| 0.5 | | 99.1 | 97.3 | 66.6 |
| 0.3 | | 99.2 | 95.8 | 73.2 |
| 0.2 | $+$ | **99.8** | **95.7** | **76.6** |
| 0.1 | | 99.3 | 93.3 | 77.4 |
| 0.0 | | 98.9 | 91.5 | 78.1 |

| $\mu$ | Fusion | $\text{Acc}_{IA}$ | $\text{Acc}_I$ | $\text{Acc}_A$ |
|---|---|---|---|---|
| 1.00 | | 98.7 | 97.5 | 57.7 |
| 0.75 | | 99.4 | 96.9 | 65.5 |
| 0.50 | | **99.7** | **95.1** | **75.1** |
| 0.25 | | 99.1 | 94.5 | 75.4 |
| 0.10 | $\oplus$ | 94.7 | 92.2 | 76.8 |
| 0.05 | | 90.3 | 89.3 | 77.7 |
| 0.01 | | 86.7 | 84.4 | 79.0 |
| 1.00 | | 99.0 | 97.3 | 66.7 |
| 0.75 | | 99.7 | 96.7 | 69.3 |
| 0.50 | | **99.8** | **95.7** | **76.6** |
| 0.25 | | 99.4 | 94.9 | 76.8 |
| 0.10 | $+$ | 97.6 | 93.9 | 78.2 |
| 0.05 | | 91.5 | 88.8 | 78.9 |
| 0.01 | | 89.6 | 83.9 | 79.4 |

value of $\mu = 1$ is equivalent to greater values of $\tau$ (e.g., $\tau = 0.5$), disabling the effects of MILES on the learning process.

**Computing u based on training metrics**. Table VI presents the results of comparing training versus validation metrics for calculating the conditional utilization rate, which is a crucial epoch-wise computation in the MILES algorithm, as outlined in Algorithm 1. The results indicate that using training metrics for conditional utilization rate computation can yield comparable outcomes. However, validation metrics generally achieve the best results more quickly during the training process compared to training set metrics. In both cases, the training metrics produced the best model at epoch 53 out of a total of 60 epochs, while the best models for validation occurred slightly earlier, at epochs 45 and 47. This demonstrates the versatility of the MILES algorithm and its potential application in tasks where a proper validation set is unavailable. We note that Table II and Table III report results of models that track the validation set metrics, for consistency and fair comparison with other methods that require the availability of a validation set (e.g., MSES).

**Recommendations**. While the impact and benefit of MILES on multimodal training may vary across datasets, fusion meth-
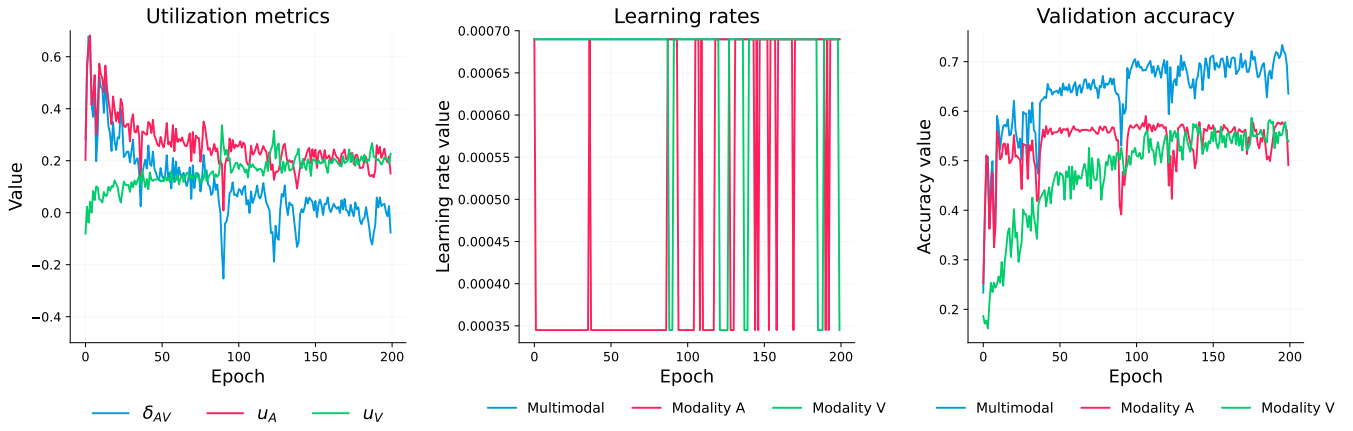
Fig. 2. Epoch-wise utilization metrics, learning rates, and validation accuracy for an example run on the CREMA-D dataset allow for the visualization of multimodal and unimodal learning dynamics using MILES.

TABLE VI
COMPARISON OF MILES RESULTS USING TRAINING VS. VALIDATION SET METRICS FOR CONDITIONAL UTILIZATION RATE CALCULATION. THE BEST RESULTS ARE SHOWN IN BOLD.

| Fusion | u | $Acc_{IA}$ | $Acc_A$ | $Acc_I$ | Epoch |
|--------|---|-----------|---------|---------|-------|
| $\oplus$ | Training | **99.8** | **95.6** | **73.6** | 53 |
| | Validation | 99.7 | 95.1 | 75.1 | **47** |
| $+$ | Training | **99.8** | **97.3** | 76.4 | 53 |
| | Validation | **99.8** | 95.7 | **76.6** | **45** |

ods, architectures, and input modalities, as shown in Table II and Table III, our empirical experience and ablation study results, enable us to provide a few recommendation for the effective usage of MILES that may work well across settings:

- Training of unimodal models can provide a good reference on the starting point for $\tau$ and $\mu$ hyper-parameters. For example, if unimodal models are very close in performance, the multimodal model may be able to integrate both modalities naturally in a more balanced way. Thus, $\tau$ could likely be set to smaller values (e.g., $\tau = 0.0$), whereas if the performance gap is significant, setting $\tau = 0.0$ may not lead to good outcomes.
- In general, start experiments with a moderate value like $\tau = 0.2$ and increase or decrease according to results.
- Initially, set $\mu$ based on the ratio between the best learning rates of unimodal models. Then, adjust $\mu$ by decreasing or increasing its value as needed. Smaller values of $\mu$ work better, as shown in Table V, because they do not excessively hinder the learning of the dominant modality.
- Once a sensitive range of values is identified for both hyper-parameters, fine-grained hyper-parameter tuning is recommended. However, in many cases, it may be sufficient to keep them fixed and vary only the learning rate.

## VI. LEARNING DYNAMICS

To visualize the multimodal and unimodal learning dynamics during training using MILES, Fig. 2 shows the epoch-wise values of utilization metrics, learning rates, and validation accuracy for an example training run on the CREMA-D dataset. Note that this is not the best run; it is provided solely for the purpose of visualizing the learning dynamics using MILES, as defined in Algorithm 1. In this example, the hyper-parameters were: (i) learning rate = 0.00068983, (ii) $\tau = 0.05$, and (iii) $\mu = 0.5$.

As can be seen in Fig. 2, the utilization of the weaker modality improves during training, with the value of $\delta_{AV} \approx 0$ (equal utilization of the modalities) at the end of the training. Throughout the training process, the learning rates are dynamically adjusted every epoch following Algorithm 1 and based on the pre-specified $\tau$ and $\mu$ values. Only one of the modality-specific learning rates is adjusted, while the multimodal learning rate is kept fixed. This allows both the utilization metric and validation accuracy to improve over time until the end of the training process.

## VII. LIMITATIONS

Seminal work on multimodal learning focus on the bimodal scenario of multimodal learning [7], [9], [10]. Similarly, all state-of-the-art baseline approaches considered in our work use two modalities to show and demonstrate their improvements [11]–[13]. Following this standard, we showcase MILES for the bimodal case, which is the most frequent setting in current state-of-the-art research related to multimodal learning [23]. In addition, most benchmark datasets are bimodal. However, MILES could be effectively adapted to scenarios with more than two modalities. For instance, for a third modality C, Algorithm 1 would need to incorporate the computation of the conditional utilization rate of the third modality ($\mathbf{u_C}$), calculate pairwise differences for all modalities ($\delta_{AB}$, $\delta_{BC}$, $\delta_{AC}$) and include a conditional block (lines 8-21 in Algorithm 1, adaptation may be required) for each difference, maintaining the same assumptions. Notwithstanding that, future work will focus on assessing the validity of these assumptions for more than two modalities and their interplay.

During our experiments, we noticed that, depending on the initialization of the network parameters, MILES might be excessively penalizing the learning of the dominant modality. This may hinder its effective learning and result in sub-optimal performance models for the dominant modality, thus requiring additional rounds of hyper-parameter tuning to achieve optimal learning. We aim to address this caveat in future work to make the proposed framework more efficient.

## VIII. Discussion and Future Work

The potential of multimodal machine learning to enhance unimodal performance is often hindered by training challenges such as modality overfitting. In this work, we propose MILES, a learning rate scheduler designed to balance and enhance multimodal machine learning. MILES leverages the epoch-wise conditional utilization rate (using validation or training performance metrics) to balance multimodal learning during training, adjusting the learning speed of the modalities to avoid both overfitting of the dominant modality and promote utilization of the non-dominant modality. Our results show that MILES effectively addresses this challenge, outperforming seven state-of-the-art baselines across datasets, tasks, fusion methods and modalities, improving the capabilities of vanilla multimodal fusion architectures for both multimodal and unimodal predictions. In addition, MILES training produces strong unimodal encoders, enabling the use of their predictive capabilities when dealing with samples that have missing modalities. MILES is governed by two hyper-parameters, which can be tuned to effectively emphasize modality learning. We provide general recommendations on how to tune these hyper-parameters to effectively balance multimodal learning during training.

Most multimodal machine learning methods and datasets focus primarily on the bimodal setting. However, as the field rapidly advances, datasets incorporating more modalities are slowly becoming available. Future work will explore extending our approach to these more complex scenarios and how the current assumptions expand to those cases, investigating the interplay between multiple modalities, and how to effectively address the training challenges of these networks.

## IX. Acknowledgment

## References

[1] S. C. Levinson and J. Holler, "The origin of human multi-modal communication," *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 369, no. 1651, p. 20130302, 2014.

[2] N. Ward, E. Paul, P. Watson, G. Cooke, C. Hillman, N. J. Cohen, A. F. Kramer, and A. K. Barbey, "Enhanced learning through multimodal training: evidence from a comprehensive cognitive, physical fitness, and neuroscience intervention," *Scientific reports*, vol. 7, no. 1, p. 5808, 2017.

[3] M. Galeotti and F. Schiantarelli, "Stock market volatility and investment: Do only fundamentals matter?" *Economica*, pp. 147–165, 1994.

[4] J. P. Langlois, "Making a diagnosis," in *Fundamentals of clinical practice*. Springer, 2002, pp. 197–217.

[5] S. Poria, E. Cambria, R. Bajpai, and A. Hussain, "A review of affective computing: From unimodal analysis to multimodal fusion," *Information fusion*, vol. 37, pp. 98–125, 2017.

[6] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018.

[7] N. Wu, S. Jastrzebski, K. Cho, and K. J. Geras, "Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks," in *International Conference on Machine Learning*. PMLR, 2022, pp. 24 043–24 055.

[8] S.-C. Huang, A. Pareek, S. Seyyedi, I. Banerjee, and M. P. Lungren, "Fusion of medical imaging and electronic health records using deep learning: a systematic review and implementation guidelines," *NPJ digital medicine*, vol. 3, no. 1, p. 136, 2020.

[9] Y. Huang, J. Lin, C. Zhou, H. Yang, and L. Huang, "Modality competition: What makes joint training of multi-modal network fail in deep learning?(provably)," in *International conference on machine learning*. PMLR, 2022, pp. 9226–9259.

[10] W. Wang, D. Tran, and M. Feiszli, "What makes training multi-modal classification networks hard?" in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 12 695–12 705.

[11] N. Fujimori, R. Endo, Y. Kawai, and T. Mochizuki, "Modality-specific learning rate control for multimodal classification," in *Pattern Recognition: 5th Asian Conference, ACPR 2019, Auckland, New Zealand, November 26–29, 2019, Revised Selected Papers, Part II 5*. Springer, 2020, pp. 412–422.

[12] Y. Yao and R. Mihalcea, "Modality-specific learning rates for effective multimodal additive late-fusion," in *Findings of the Association for Computational Linguistics: ACL 2022*, 2022, pp. 1824–1834.

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

[14] A. Guerra-Manzanares and F. Shamout, "MIND: Modality-informed knowledge distillation framework for multimodal clinical prediction tasks," *Transactions on Machine Learning Research*, 2025. [Online]. Available: https://openreview.net/forum?id=BhOJreYmur

[15] H. Cao, D. G. Cooper, M. K. Keutmann, R. C. Gur, A. Nenkova, and R. Verma, "Crema-d: Crowd-sourced emotional multimodal actors dataset," *IEEE transactions on affective computing*, vol. 5, no. 4, pp. 377–390, 2014.

[16] L. Khacef, L. Rodriguez, and B. Miramond, "Written and spoken digits database for multimodal learning," https://zenodo.org/doi/10.5281/zenodo.3515934, 2019.

[17] Y. LeCun, "The mnist database of handwritten digits," *http://yann. lecun. com/exdb/mnist/*, 1998.

[18] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.

[19] G. Bezirganyan, S. Sellami, L. Berti-Équille, and S. Fournier, "Luma: A benchmark dataset for learning from uncertain and multimodal data," *arXiv preprint arXiv:2406.09864*, 2024.

[20] J. Arevalo, T. Solorio, M. Montes-y Gómez, and F. A. González, "Gated multimodal units for information fusion," *arXiv preprint arXiv:1702.01992*, 2017.

[21] J. D. M.-W. C. Kenton and L. K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of naacL-HLT*, vol. 1. Minneapolis, Minnesota, 2019, p. 2.

[22] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[23] P. Xu, X. Zhu, and D. A. Clifton, "Multimodal learning with transformers: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 10, pp. 12 113–12 132, 2023.