

# Toward Robust Incomplete Multimodal Sentiment Analysis via Hierarchical Representation Learning

Mingcheng Li<sup>1,3\*</sup> Dingkang Yang<sup>1,3\*†</sup> Yang Liu<sup>1</sup> Shunli Wang<sup>1,3</sup> Jiawei Chen<sup>1,3</sup>  
 Shuaibing Wang<sup>1,3</sup> Jinjie Wei<sup>1,3</sup> Yue Jiang<sup>1,3</sup> Qingyao Xu<sup>1,3</sup> Xiaolu Hou<sup>1,3</sup>  
 Mingyang Sun<sup>1,3</sup> Ziyun Qian<sup>1,3</sup> Dongliang Kou<sup>1,3</sup> Lihua Zhang<sup>1,2,3,4,5†</sup>

<sup>1</sup>Academy for Engineering and Technology, Fudan University, Shanghai, China

<sup>2</sup>Institute of Metaverse & Intelligent Medicine, Fudan University, Shanghai, China

<sup>3</sup>Cognition and Intelligent Technology Laboratory, Shanghai, China

<sup>4</sup>Jilin Provincial Key Laboratory of Intelligence Science and Engineering, Changchun, China

<sup>5</sup>Engineering Research Center of AI and Robotics, Ministry of Education, Shanghai, China.

mingchengli21@m.fudan.edu.cn, dkyang20@fudan.edu.cn

## Abstract

Multimodal Sentiment Analysis (MSA) is an important research area that aims to understand and recognize human sentiment through multiple modalities. The complementary information provided by multimodal fusion promotes better sentiment analysis compared to utilizing only a single modality. Nevertheless, in real-world applications, many unavoidable factors may lead to situations of uncertain modality missing, thus hindering the effectiveness of multimodal modeling and degrading the model’s performance. To this end, we propose a Hierarchical Representation Learning Framework (HRLF) for the MSA task under uncertain missing modalities. Specifically, we propose a fine-grained representation factorization module that sufficiently extracts valuable sentiment information by factorizing modality into sentiment-relevant and modality-specific representations through crossmodal translation and sentiment semantic reconstruction. Moreover, a hierarchical mutual information maximization mechanism is introduced to incrementally maximize the mutual information between multi-scale representations to align and reconstruct the high-level semantics in the representations. Ultimately, we propose a hierarchical adversarial learning mechanism that further aligns and adapts the latent distribution of sentiment-relevant representations to produce robust joint multimodal representations. Comprehensive experiments on three datasets demonstrate that HRLF significantly improves MSA performance under uncertain modality missing cases.

## 1 Introduction

Multimodal sentiment analysis (MSA) has attracted wide attention in recent years. Unlike unimodal emotion recognition tasks [9, 63, 64, 53, 56], MSA understands and recognizes human emotions through multiple modalities, including language, audio, and visual [31, 58]. Previous studies have shown that combining complementary information among different modalities facilitates valuable semantic generation [41, 40, 61, 55, 62]. MSA has been well studied so far under the assumption that all modalities are available in the training and inference phases [12, 66, 54, 57, 56, 25, 59, 60]. Nevertheless, in real-world applications, modalities may be missing due to security concerns, background noises, sensor limitations and so on. Ultimately, these incomplete multimodal data significantly hinder the performance of MSA. For instance, as shown in Figure 1, the entire visual

\*Equal contributions. †Corresponding authors.

modality and some frame-level features in the language and audio modalities are missing, leading to an incorrect prediction.

In recent years, many studies [8, 28, 26, 49, 37, 50, 76, 74, 68, 27, 23, 22] attempt to address the problem of missing modalities in MSA. For example, SMIL [29] estimates the latent features of the missing modality data via Bayesian Meta-Learning. However, these methods are constrained by the following factors: **(i)** Implementing complex feature interactions for incomplete modalities leads to a large amount of information redundancy and cumulative errors, resulting in ineffective extraction of sentiment semantics. **(ii)** Lacking consideration of semantic and distributional alignment of representations, causing imprecise feature reconstruction and nonrobust joint representations.

To address the above issues, we propose a Hierarchical Representation Learning Framework (HRLF) for the MSA task under uncertain missing modalities. HRLF has three core contributions: **(i)** We present a fine-grained representation factorization module that sufficiently extracts valuable sentiment information by factorizing modality into sentiment-relevant and modality-specific representations through intra- and inter-modality translations and sentiment semantic reconstruction. **(ii)** Furthermore, a hierarchical mutual information maximization mechanism is introduced to incrementally align the high-level semantics by maximizing the mutual information of the multi-scale representations of both networks in knowledge distillation. **(iii)** Eventually, we propose a hierarchical adversarial learning mechanism to progressively align the latent distributions of representations leveraging multi-scale adversarial learning. Based on these components, HRLF significantly improves MSA performance under uncertain modality missing cases on three multimodal benchmarks.

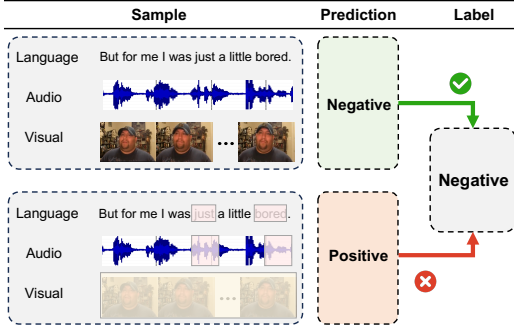


Figure 1: A case of incorrect prediction by the traditional model with missing modalities. The pink and yellow areas indicate intra- and inter-modality missingness, respectively.

## 2 Related Work

### 2.1 Multimodal Sentiment Analysis

Multimodal Sentiment Analysis (MSA) seeks to comprehend and analyze human sentiment by utilizing diverse modalities. Unlike conventional single-modality sentiment recognition, MSA poses greater challenges owing to the intricate nature of processing and analyzing heterogeneous data across modalities. Mainstream studies in MSA [69, 70, 44, 12, 11, 42, 25] focus on designing complex fusion paradigms and interaction mechanisms to improve MSA performance. For instance, CubeMLP [42] employs three distinct multi-layer perceptron units for feature amalgamation along three axes. However, these methods rely on complete modalities and thus are impractical for real-world deployment. There are two primary approaches for addressing the missing modality problem in MSA: (1) Generative methods [8, 28, 26, 49] and (2) joint learning methods [37, 50, 76, 74, 68, 27]. Generative methods aim to regenerate missing features and semantics within modalities by leveraging the distributions of available modalities. For example, TFR-Net [67] employs a feature reconstruction module to guide the extractor to reconstruct missing semantics. Joint learning methods focus on deriving cohesive joint multimodal representations based on inter-modality correlations. For instance, MMIN [76] produces robust joint multimodal representations via cross-modality imagination. However, these methods cannot extract rich sentiment information from incomplete modalities due to their inefficient interaction. In contrast, our learning paradigm achieves effective extraction and precise reconstruction of sentiment semantics through complete modality factorization.

### 2.2 Factorized Representation Learning

The fundamental goal of learning factorized representations is to disentangle representations that have different semantics and distributions. This separation enables the model to more effectively capture intrinsic information and yield favorable modality representations. Previous methods of factorized

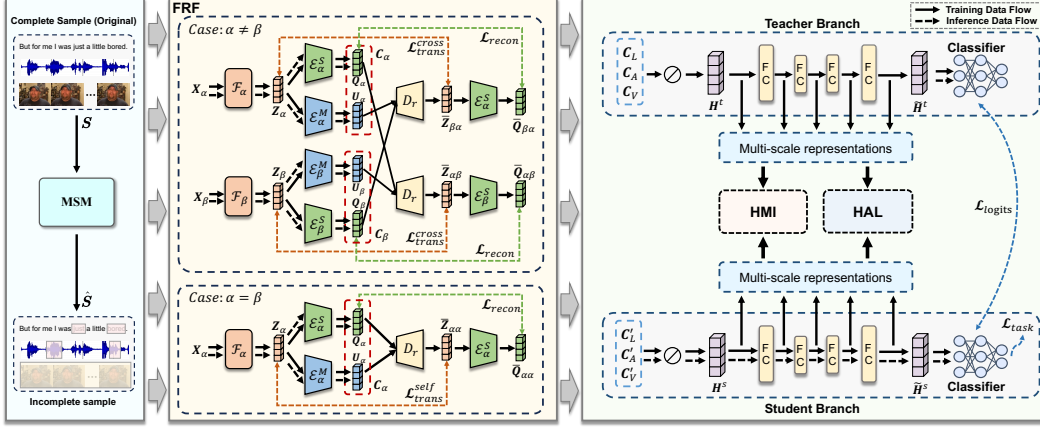


Figure 2: The structure of our HRLF, which consists of three core components: Fine-grained Representation Factorization (FRF) module, Hierarchical Mutual Information (HMI) maximization mechanism, and Hierarchical Adversarial Learning (HAL) mechanism.

representation learning primarily rely on auto-encoders [3] and generative adversarial networks [32]. For example, FactorVA [18] is introduced to achieve factorization by leveraging the characteristic that representations are both factorial and independent in dimension. Recently, factorization learning has been progressively utilized in MSA tasks [54, 25, 57]. For instance, FDMER [54] utilizes consistency and discreteness constraints between modalities to disentangle modalities into modality-invariant and modality-private features. DMD [25] disentangles each modality into modality-independent and modality-exclusive representations and then implements a knowledge distillation strategy among the representations with dynamic graphs. MFSA [57] refines multimodal representations and learns complementary representations across modalities by learning modality-specific and modality-agnostic representations. Despite the progress these studies have brought to MSA, certain limitations persist: (i) The supervision of the factorization process is coarse-grained and insufficient. (ii) Focusing solely on factorizing distinct representations at the modality level, without taking into account sentimentally beneficial and relevant representations. By contrast, the proposed method decomposes sentiment-relevant representations precisely through intra- and inter-modality translation and sentiment semantic reconstruction. Furthermore, hierarchical mutual information maximization and adversarial learning paradigms are employed to refine and optimize the representation of factorization at the semantic level and the distributional level, respectively, thus yielding robust joint multimodal representations.

### 2.3 Knowledge Distillation

Knowledge distillation leverages additional supervisory signals from a pre-trained teacher network to aid in training a student network [15]. There are generally two categories of knowledge distillation methods: distillation from intermediate features [13, 14, 19, 33, 35, 39, 45, 43, 65, 73] and distillation from logits [6, 10, 30, 52, 75]. Many studies [5, 17, 38, 21, 47, 51] utilize knowledge distillation for MSA tasks with missing modalities. These approaches aim to transfer dark knowledge from teacher networks trained on complete modalities to student networks trained by missing modalities. The teacher network typically provides richer and more comprehensive feature representations than the student network. For instance, KD-Net [17] utilizes a teacher network with complete modalities to supervise the unimodal student network at both the feature and logits levels. Despite their promising results, these methods neglect precise supervision of representations, resulting in low-quality knowledge transfer. To this end, we implement hierarchical semantic and distributional alignment of the multi-scale representations of both networks to transfer knowledge effectively.

## 3 Methodology

### 3.1 Problem Formulation

Given a multimodal video segment with three modalities as  $\mathcal{S} = [\mathbf{X}_L, \mathbf{X}_A, \mathbf{X}_V]$ , where  $\mathbf{X}_L \in \mathbb{R}^{T_L \times d_L}$ ,  $\mathbf{X}_A \in \mathbb{R}^{T_A \times d_A}$ , and  $\mathbf{X}_V \in \mathbb{R}^{T_V \times d_V}$  denote language, audio, and visual modalities,

respectively.  $\mu = \{L, A, V\}$  denotes the set of modality types.  $T_m(\cdot)$  is the sequence length and  $d_m(\cdot)$  is the embedding dimension, where  $m \in \mu$ . We define two missing modality cases to simulate the most natural and holistic challenges in real-world scenarios: (1) *intra-modality missingness*, which indicates some frame-level features in the modality sequences are missing. (2) *inter-modality missingness*, which denotes some modalities are entirely missing. We aim to recognize the utterance-level sentiments using incomplete multimodal data.

### 3.2 Overall Framework

Figure 2 illustrates the main workflow of HRLF. The teacher and student networks adopt a consistent structure but have different parameters. During the training phase, the workflow of our HRLF is as follows: (i) We first train the teacher network with complete-modality samples and their sentiment labels. (ii) Given a video segment sample  $\mathcal{S}$ , we generate a missing-modality sample  $\tilde{\mathcal{S}}$  with the Modality Stochastic Missing (MSM) strategy. MSM simultaneously performs intra-modality missingness and inter-modality missingness.  $\mathcal{S}$  and  $\tilde{\mathcal{S}}$  are fed into the pre-trained teacher network and the initialized student network, respectively. (iii) We input each sample into the FRF module, to factorize each modality into a sentiment-relevant representation  $\mathbf{Q}_m$  and a modality-specific representation  $\mathbf{U}_m$ , where  $m \in \mu$ . (iv) Sequences  $[\mathbf{C}_L, \mathbf{C}_A, \mathbf{C}_V]$  and  $[\mathbf{C}'_L, \mathbf{C}'_A, \mathbf{C}'_V]$  are generated by concatenating  $\mathbf{Q}_m$  and  $\mathbf{U}_m$  from all modalities in the teacher and student networks. Each element of the sequences is concatenated to yield the joint multimodal representations  $\mathbf{H}^t$  and  $\mathbf{H}^s$ . (v) The multi-scale representations of both networks are obtained by passing  $\mathbf{H}^t$  and  $\mathbf{H}^s$  through the fully-connected layers. The proposed HMI and HAL are used to align the semantics and distribution between the multiscale representations. (vi) The outputs  $\tilde{\mathbf{H}}^t$  and  $\tilde{\mathbf{H}}^s$  of the fully-connected layers are fed into the task-specific classifier to get logits  $\mathbf{L}^t$  and  $\mathbf{L}^s$ . We constrain the consistency between logits and utilize  $\mathbf{L}^s$  to implement the sentiment prediction. In the inference phase, testing samples are only fed into the student network for downstream tasks.

### 3.3 Fine-grained Representation Factorization

Modality missing leads to ambiguous sentiment cues in the modality and information redundancy in multimodal fusion. It hinders the model from capturing valuable sentiment semantics and filtering sentiment irrelevant information. Although previous studies in MSA [12, 54] decompose the task-relevant semantics contained in the modality to some extent via simple auto-encoder networks with reconstruction constraints, their purification of sentiment semantics is inadequate, and they cannot be applied to modality missing scenarios. Therefore, we propose a Fine-grained Representation Factorization (FRF) module to capture sentiment semantics in modalities. The core idea is to factorize each modality representation into two types of representations: (1) sentiment-relevant representation, which contains the holistic sentiment semantics of the sample. It is modality-independent, shared across all modalities of the same subject, and robust to modality missing situations. (2) modality-specific representation, which represents modality-specific task-independent information.

As shown in Figure 2, FRF receives the multimodal sequences  $[\mathbf{X}_L, \mathbf{X}_A, \mathbf{X}_V]$  with modality number  $n = 3$ . The modality  $\mathbf{X}_\alpha$  with  $\alpha \in \mu$  passes through a 1D temporal convolutional layer with kernel size  $3 \times 3$  and adds the positional embedding [46] to obtain the preliminary representations, denoted as  $\mathbf{R}_\alpha = \mathbf{W}_{3 \times 3}(\mathbf{X}_\alpha) + PE(T_\alpha, d) \in \mathbb{R}^{T_\alpha \times d}$ . The  $\mathbf{R}_\alpha$  is fed into a Transformer [46] encoder  $\mathcal{F}_\alpha(\cdot)$ , and the last element of its output is denoted as  $\mathbf{Z}_\alpha = \mathcal{F}_\alpha(\mathbf{R}_\alpha) \in \mathbb{R}^d$ . The  $\mathbf{Z}_\alpha \in \mathcal{Z}_\alpha$  is the low-level modality representation of the modality  $\alpha$ . We aim to factorize modality representation  $\mathbf{Z}_\alpha$  into a sentiment-relevant representation  $\mathbf{Q}_\alpha$  by a sentiment encoder  $\mathbf{Q}_\alpha = \mathcal{E}_\alpha^S(\mathbf{Z}_\alpha)$  and a modality-specific representation  $\mathbf{U}_\alpha$  by a modality encoder  $\mathbf{U}_\alpha = \mathcal{E}_\alpha^M(\mathbf{Z}_\alpha)$ .  $\mathcal{E}_\alpha^S(\cdot)$  and  $\mathcal{E}_\alpha^M(\cdot)$  are composed of multi-layer perceptrons with the ReLU activation. The following two processes ensure adequate factorization and semantic reinforcement of the above two representations.

**Intra- and Inter-modality Translation.** The proposed FRF effectively decouples sentiment-relevant and modality-specific representations by simultaneously performing intra- and inter-modality translations. Given a pair of representations  $\mathbf{Q}_\alpha$  and  $\mathbf{U}_\beta$  factorized by  $\mathbf{Z}_\alpha$  and  $\mathbf{Z}_\beta$  with  $\alpha, \beta \in \mu$ , the decoder  $\mathcal{D}_r(\cdot)$  is supposed to translate and synthesize the representation  $\bar{\mathbf{Z}}_{\alpha\beta}$ , whose reconstructed domain corresponds to the modality representation  $\mathbf{Z}_\beta \in \mathcal{Z}_\beta$ . The  $\mathcal{D}_r(\cdot)$  consists of feed-forward neural layers. The modality translations include intra-modality translation (*i.e.*,  $\alpha = \beta$ ) and inter-modality translation (*i.e.*,  $\alpha \neq \beta$ ), whose losses are respectively denoted as:

$$\mathcal{L}_{trans}^{self} = \frac{1}{n} \sum_{\alpha \in \mu} \mathbf{E}_{\mathbf{Z}_\alpha \sim \mathcal{Z}_\alpha} [\|\bar{\mathbf{Z}}_{\alpha\alpha} - \mathbf{Z}_\alpha\|_2], \quad (1)$$

$$\mathcal{L}_{trans}^{cross} = \frac{1}{n^2 - n} \sum_{\alpha \in \mu} \sum_{\beta \in \mu, \beta \neq \alpha} \mathbf{E}_{\mathbf{Z}_\alpha \sim \mathcal{Z}_\alpha, \mathbf{Z}_\beta \sim \mathcal{Z}_\beta} [\|\bar{\mathbf{Z}}_{\alpha\beta} - \mathbf{Z}_\beta\|_2], \quad (2)$$

where  $\bar{\mathbf{Z}}_{\alpha\beta} = \mathcal{D}_r(\mathcal{E}_\alpha^S(\mathbf{Z}_\alpha), \mathcal{E}_\beta^M(\mathbf{Z}_\beta))$ . The translation loss is denoted as:  $\mathcal{L}_{trans} = \mathcal{L}_{trans}^{self} + \mathcal{L}_{trans}^{cross}$ .

**Sentiment Semantic Reconstruction.** To ensure that the reconstructed modality still contains the sentiment semantics from the original modality, we encourage both to maintain the consistency of sentiment-relevant semantics and utilize the following loss for constraints, denoted as:

$$\mathcal{L}_{recon} = \frac{1}{n^2} \sum_{\alpha \in \mu} \sum_{\beta \in \mu} \mathbf{E}_{\mathbf{Z}_\alpha \sim \mathcal{Z}_\alpha, \mathbf{Z}_\beta \sim \mathcal{Z}_\beta} [\|\bar{\mathbf{Q}}_{\beta\alpha} - \mathbf{Q}_\alpha\|_2], \quad (3)$$

where  $\bar{\mathbf{Q}}_{\beta\alpha} = \mathcal{E}_\alpha^S(\mathcal{D}_r(\mathcal{E}_\beta^S(\mathbf{Z}_\beta), \mathcal{E}_\alpha^M(\mathbf{Z}_\alpha)))$  is the sentiment-relevant representation derived from the reconstructed modality representation. Consequently, the final loss of the FRF is denoted as:

$$\mathcal{L}_{FRF} = \mathcal{L}_{trans} + \mathcal{L}_{recon}. \quad (4)$$

### 3.4 Hierarchical Mutual Information Maximization

The underlying assumption of knowledge distillation is that layers in the pre-trained teacher network can represent certain attributes of given inputs that exist in the task [15]. For successful knowledge transfer, the student network must learn to incorporate such attributes into its own learning. Nevertheless, previous studies [17, 38, 21] based on knowledge distillation simply constrain the consistency between the features of both networks and lack consideration of the intrinsic semantics and inherent properties of the features, leading to semantic misalignment. From the perspective of information theory [1], semantic alignment and attribute mining of representations can be characterized as maintaining high mutual information among the layers of the teacher and student networks. We construct a Hierarchical Mutual Information (HMI) maximization mechanism to implement sufficient semantic alignment and maximize mutual information. The core idea is to progressively align the semantics of representations through a hierarchical learning paradigm.

Specifically, the sentiment-relevant and modality-specific representations  $\mathbf{Q}_m$  and  $\mathbf{U}_m$  of all modalities for teacher and student networks are concatenated to obtain the sequences  $[\mathbf{C}_L, \mathbf{C}_A, \mathbf{C}_V]$  and  $[\mathbf{C}'_L, \mathbf{C}'_A, \mathbf{C}'_V]$ . Each element of the sequences is concatenated to yield the joint multimodal representations  $\mathbf{H}^t$  and  $\mathbf{H}^s$ . The fully-connected layers are utilized to refine the representation  $\mathbf{H}^w \in \mathbb{R}^{3d}$  with  $w \in \{t, s\}$ , yielding  $\bar{\mathbf{H}}^w \in \mathbb{R}^{3d}$ . Moreover, we obtain the intermediate multi-scale representations of all layers, denoted as  $\mathbf{I}_1^w \in \mathbb{R}^{2d}$ ,  $\mathbf{I}_2^w \in \mathbb{R}^d$ , and  $\mathbf{I}_3^w \in \mathbb{R}^{2d}$ . For the above five representations, we concatenate features of the same scale to obtain multi-scale representations  $\mathbf{E}_1^w \in \mathbb{R}^{3d}$ ,  $\mathbf{E}_2^w \in \mathbb{R}^{2d}$ , and  $\mathbf{E}_3^w \in \mathbb{R}^d$ , which are utilized in the subsequent computation.

To estimate and compute the mutual information between representations, we define two random variables  $\mathbf{X}$  and  $\mathbf{Y}$ . The  $P(\mathbf{X})$  and  $P(\mathbf{Y})$  are the marginal probability density function of  $\mathbf{X}$  and  $\mathbf{Y}$ . The joint probability density function of  $\mathbf{X}$  and  $\mathbf{Y}$  is denoted as  $P(\mathbf{X}, \mathbf{Y})$ . The mutual information of the random variables  $\mathbf{X}$  and  $\mathbf{Y}$  is represented as:

$$I(\mathbf{X}; \mathbf{Y}) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} \right]. \quad (5)$$

We only need to obtain the maximum value of the mutual information, without focusing on its exact value. Referring to Deep InfoMax [16], we estimate the mutual information between variables based on the Jensen-Shannon Divergence (JSD). The mutual information maximization issue translates into minimizing the JSD between the joint distribution  $p(\mathbf{x}, \mathbf{y})$  and the marginal distribution  $p(\mathbf{x})p(\mathbf{y})$ .

$$JSD(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) = \frac{1}{2} (D_{KL}(p(\mathbf{x}, \mathbf{y}) \| m) + D_{KL}(p(\mathbf{x})p(\mathbf{y}) \| m)), \quad (6)$$

where  $m = \frac{1}{2}(p(\mathbf{x}, \mathbf{y}) + p(\mathbf{x})p(\mathbf{y}))$  and  $D_{KL}$  is Kullback-Leibler divergence. Mutual information maximization is achieved by maximizing the dyadic lower bound of JSD, denoted as:

$$MI(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{P(\mathbf{x}, \mathbf{y})} [-sp(-\mathcal{T}_\theta(\mathbf{x}, \mathbf{y}))] + \mathbb{E}_{P(\mathbf{x})P(\mathbf{y})} [-sp(\mathcal{T}_\theta(\mathbf{x}, \mathbf{y}))], \quad (7)$$

where  $sp(w) = \log(1 + e^w)$  and  $\mathcal{T}_\theta(x, y)$  is the statistics network which is a neural network with parameters  $\theta$ . HMI maximizes the mutual information between hierarchical representations in knowledge distillation, whose optimization loss is expressed as:

$$\mathcal{L}_{HMI} = - \sum_{i=1}^3 MI(\mathbf{E}_i^t, \mathbf{E}_i^s). \quad (8)$$

### 3.5 Hierarchical Adversarial Learning

Considering that the teacher network has more robust and stable representation distributions, we also need to encourage the alignment of representation distributions in the latent space. Traditional methods [38, 17, 21] simply minimize the KL divergence between both networks, which easily disturbs the underlying learning of the student network in the deep layers, leading to confounded distributions and unrobust joint multimodal representations.

To this end, we propose a Hierarchical Adversarial Learning (HAL) mechanism for incrementally aligning the latent distributions between representations of student and teacher networks. The central principle is that the student network tries to generate representations to mislead the discriminator  $\mathcal{D}_e(\cdot)$ , while  $\mathcal{D}_e(\cdot)$  discriminates between the representations of the student and teacher networks. In practice,  $\mathcal{D}_e(\cdot)$  is the fully-connected layers. Specifically, given multi-scale representations of  $\mathbf{E}_1^w \in \mathbb{R}^{3d}$ ,  $\mathbf{E}_2^w \in \mathbb{R}^{2d}$ , and  $\mathbf{E}_3^w \in \mathbb{R}^d$  with  $w \in \{t, s\}$ , we implement adversarial learning on the same-scale representations of the teacher and student networks to hierarchically supervise consistency. The objective function of HAL is formatted as:

$$\mathcal{L}_{HAL} = \sum_{i=1}^3 \log(1 - \mathcal{D}_e(\mathbf{E}_i^s)) + \log(\mathcal{D}_e(\mathbf{E}_i^t)). \quad (9)$$

### 3.6 Optimization Objectives

The  $\tilde{H}^t$  and  $\tilde{H}^s$  of the teacher and student networks are fed into their task-specific classifiers to produce logits  $\mathbf{L}^t$  and  $\mathbf{L}^s$ , respectively, and the consistency of both is constrained with KL divergence loss, denoted as  $\mathcal{L}_{KL} = KL(\mathbf{L}^t, \mathbf{L}^s)$ . The  $\mathbf{L}^s$  is used for sentiment recognition and supervised with task loss, represented as  $\mathcal{L}_{task}$ . For the classification and regression tasks, we use cross-entropy and MSE loss as the task losses, respectively. The overall training objective  $\mathcal{L}_{total}$  is expressed as  $\mathcal{L}_{total} = \mathcal{L}_{task} + \mathcal{L}_{FRF} + \mathcal{L}_{HMI} + \mathcal{L}_{HAL} + \mathcal{L}_{KL}$ .

## 4 Experiments

### 4.1 Datasets and Evaluation Metrics

We conduct our experiments on three MSA benchmarks, including MOSI [71], MOSEI [72], and IEMOCAP [4]. The experiments are performed under the word-aligned setting. MOSI is a realistic dataset for MSA. It comprises 2,199 short monologue video clips taken from 93 YouTube movie review videos. There are 1,284, 229, and 686 video clips in train, valid, and test data, respectively. MOSEI is a dataset consisting of 22,856 movie review video clips, which has 16,326, 1,871, and 4,659 samples in train, valid, and test data. Each sample of MOSI and MOSEI is labelled by human annotators with a sentiment score of -3 (strongly negative) to +3 (strongly positive). On the MOSI and MOSEI datasets, we utilize two evaluation metrics, including the Mean Absolute Error (MAE) and F1 score computed for positive/negative classification results. The IEMOCAP dataset consists of 4,453 samples of video clips. Its predetermined data partition has 2,717, 798, and 938 samples in train, valid, and test data. As recommended by [48], four emotions (*i.e.*, happy, sad, angry, and neutral) are selected for emotion recognition. The F1 score is used as the metric.

### 4.2 Implementation Details

**Feature Extraction.** The Glove embedding [36] is used to convert the video transcripts to obtain a 300-dimensional vector for the language modality. For the audio modality, we employ the COVAREP

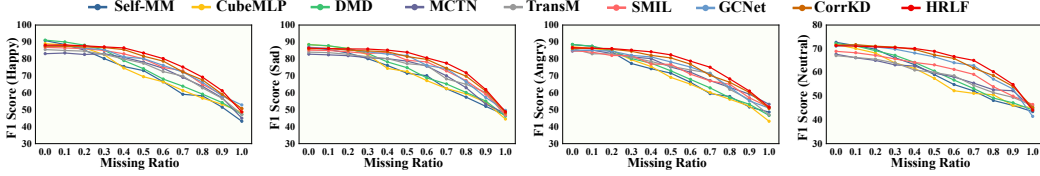


Figure 3: Comparison results of intra-modality missingness on IEMOCAP. We report on the F1 score metric for the happy, sad, angry, and neutral categories.

toolkit [7] to extract 74-dimensional acoustic features, including 12 Mel-frequency cepstral coefficients, voiced/unvoiced segmenting features, and glottal source parameters. For the visual modality, we utilize the Facet [2] to indicate 35 facial action units that record facial movement.

**Experimental Setup.** Regarding the MOSI [71] and MOSEI [72] datasets, we use the aligned multimodal sequences therein (*e.g.*, all sequences of modalities have length 300) as the original input for the HRLF. All models are built on the Pytorch [34] toolbox with four NVIDIA Tesla V100 GPUs. The Adam optimizer [20] is employed for network optimization. For MOSI, MOSEI, and IEMOCAP, the detailed hyper-parameter settings are as follows: the learning rates are  $\{1e-3, 2e-3, 4e-3\}$ , the batch sizes are  $\{128, 16, 32\}$ , the epoch numbers are  $\{50, 20, 30\}$ , and the attention heads are  $\{10, 8, 10\}$ . The embedding dimension is 40 on all three datasets. The raw features at the modality missing positions are replaced by zero vectors. For a fair comparison, we re-implement the State-Of-The-Art (SOTA) methods and combine them with our experimental paradigms. All experimental results are averaged over multiple experiments using five different random seeds.

### 4.3 Comparison with State-of-the-art Methods

We conduct a comparison between HRLF and eight representative, reproducible state-of-the-art (SOTA) methods, including complete-modality methods: Self-MM [66], CubeMLP [42], and DMD [25], and missing-modality methods: 1) joint learning methods (*i.e.*, MCTN [37], TransM [50], and CorrKD [24]), and 2) generative methods (*i.e.*, SMIL [29] and GCNet [26]). The extensive experiments are designed to comprehensively assess the robustness and effectiveness of HRLF in scenarios involving both intra-modality and inter-modality missingness.

**Robustness to Intra-modality Missingness.** We simulate intra-modality missingness by randomly discarding frame-level features in sequences with ratio  $p \in \{0.1, 0.2, \dots, 1.0\}$ . To visualize the robustness of all models, Figure 3 and 4 show the performance curves of the models for different ratios  $p$ . We have the following important observations. (i) As the ratio  $p$  increases, the performance of all models declines. This phenomenon demonstrates that intra-modality missingness leads to significant sentiment semantic loss and fragile multimodal representations. (ii) Compared to complete-modality methods, our HRLF demonstrates notable performance advantages in missing-modality testing conditions and competitive performance in complete-modality testing conditions. This is because complete-modality methods rely on the assumption of data completeness, while training paradigms for missing modalities excel in capturing and reconstructing valuable sentiment semantics from incomplete multimodal data. (iii) In contrast to the missing-modality methods, our HRLF demonstrates the highest level of robustness. Through the purification of sentiment semantics and the dual alignment of representations, the student network masters the core competencies of precisely reconstructing missing semantics and generating robust multimodal representations.

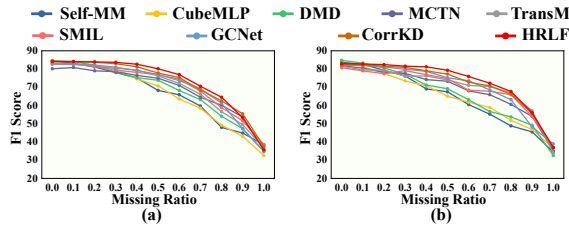


Figure 4: Comparison results of intra-modality missingness on (a) MOSI and (b) MOSEI.

**Robustness to Inter-modality Missingness.** To simulate the case of inter-modality missingness, we remove certain entire modalities from the samples. Tables 1 and 2 contrast the models’ resilience to inter-modality missingness. The notation “ $\{l\}$ ” signifies that only the language modality is available,

Table 1: Comparison of performance under six possible testing conditions of inter-modality missingness and the complete-modality testing condition on the MOSI and MOSEI datasets. T-test is conducted on ‘‘Avg.’’ column. \* indicates that  $p < 0.05$  (compared with the SOTA CorrKD).

Datasets	Models	Testing Conditions							
		{l}	{a}	{v}	{l, a}	{l, v}	{a, v}	Avg.	{l, a, v}
MOSI	Self-MM [66]	67.80	40.95	38.52	69.81	74.97	47.12	56.53	<b>84.64</b>
	CubeMLP [42]	64.15	38.91	43.24	63.76	65.12	47.92	53.85	84.57
	DMD [25]	68.97	43.33	42.26	70.51	68.45	50.47	57.33	84.50
	MCTN [37]	75.21	59.25	58.57	77.81	74.82	64.21	68.31	80.12
	TransM [50]	77.64	63.57	56.48	82.07	80.90	67.24	71.32	82.57
	SMIL [29]	78.26	67.69	59.67	79.82	79.15	71.24	72.64	82.85
	GCNet [26]	80.91	65.07	58.70	<b>84.73</b>	<b>83.58</b>	70.02	73.84	83.20
	CorrKD [24]	81.20	66.52	60.72	83.56	82.41	73.74	74.69	83.94
<b>HRLF (Ours)</b>	<b>83.36</b>	<b>69.47</b>	<b>64.59</b>	83.82	83.56	<b>75.62</b>	<b>76.74*</b>	84.15	
MOSEI	Self-MM [66]	71.53	43.57	37.61	75.91	74.62	49.52	58.79	83.69
	CubeMLP [42]	67.52	39.54	32.58	71.69	70.06	48.54	54.99	83.17
	DMD [25]	70.26	46.18	39.84	74.78	72.45	52.70	59.37	<b>84.78</b>
	MCTN [37]	75.50	62.72	59.46	76.64	77.13	64.84	69.38	81.75
	TransM [50]	77.98	63.68	58.67	80.46	78.61	62.24	70.27	81.48
	SMIL [29]	76.57	65.96	60.57	77.68	76.24	66.87	70.65	80.74
	GCNet [26]	80.52	66.54	61.83	81.96	81.15	69.21	73.54	82.35
	CorrKD [24]	80.76	66.09	62.30	81.74	<b>81.28</b>	71.92	74.02	82.16
<b>HRLF (Ours)</b>	<b>82.05</b>	<b>69.32</b>	<b>64.90</b>	<b>82.62</b>	81.09	<b>73.80</b>	<b>75.63*</b>	82.93	

while the audio and visual modalities are missing. ‘‘{l, a, v}’’ denotes the complete-modality testing condition where all modalities are available. ‘‘Avg.’’ indicates the average performance across six missing-modality testing conditions. We have the following key findings: **(i)** The inter-modality missingness leads to a decline in performance for all models, indicating that integrating complementary information from diverse modalities enhances the sentiment semantics within joint representations. **(ii)** Across all six testing conditions involving inter-modality missingness, our HRLF consistently demonstrates superior performance among the majority of metrics, affirming its robustness. For example, on the MOSI dataset, HRLF’s average F1 score is improved by 2.05% compared to CorrKD, and in particular by 3.87% in the testing condition where only visual modality is available (*i.e.*, {v}). The advantage comes from its learning of fine-grained representation factorization and the hierarchical semantic alignment and distributional alignment. **(iii)** In unimodal testing scenarios, HRLF’s performance using only the language modality significantly exceeds other configurations, showing performance similar to that of the complete-modality setup. In bimodal testing scenarios, configurations involving the language modality exhibit superior performance, even outperforming the complete-modality setup in specific metrics. This phenomenon underscores the richness of sentiment semantics within the language modality and its dominance in sentiment inference and missing semantic reconstruction processes.

#### 4.4 Ablation Studies

To affirm the effectiveness and indispensability of the module and mechanisms and strategies proposed in HRLF, we perform ablation experiments under two missing-modality scenarios on the MOSI dataset, as shown in Table 3 and Figure 5. We have the following important observations: **(i)** First, when the FRF is removed, sentiment-relevant and modality-specific information in the modalities are confused, hindering sentiment recognition and leading to significant performance degradation. This phenomenon demonstrates the effectiveness of the proposed representation factorization paradigm for adequate capture of valuable sentiment semantics. **(ii)** When our HMI

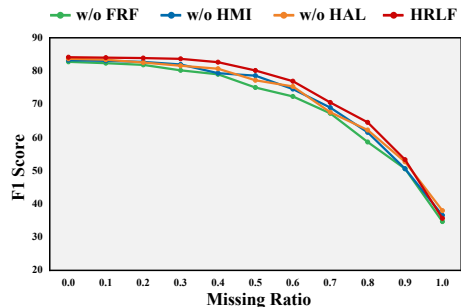


Figure 5: Ablation results of intra-modality missingness case on the MOSI dataset.



Table 2: Comparison of performance under six possible testing conditions of inter-modality missingness and the complete-modality testing condition on the IEMOCAP dataset. T-test is conducted on ‘‘Avg.’’ column. \* indicates that  $p < 0.05$  (compared with the SOTA CorrKD).

Models	Metrics	Testing Conditions							
		{l}	{a}	{v}	{l, a}	{l, v}	{a, v}	Avg.	{l, a, v}
Self-MM [66]	Happy	66.9	52.2	50.1	69.9	68.3	56.3	60.6	90.8
	Sad	68.7	51.9	54.8	71.3	69.5	57.5	62.3	86.7
	Angry	65.4	53.0	51.9	69.5	67.7	56.6	60.7	88.4
	Neutral	55.8	48.2	50.4	58.1	56.5	52.8	53.6	<b>72.7</b>
CubeMLP [42]	Happy	68.9	54.3	51.4	72.1	69.8	60.6	62.9	89.0
	Sad	65.3	54.8	53.2	70.3	68.7	58.1	61.7	<b>88.5</b>
	Angry	65.8	53.1	50.4	69.5	69.0	54.8	60.4	87.2
	Neutral	53.5	50.8	48.7	57.3	54.5	51.8	52.8	71.8
DMD [25]	Happy	69.5	55.4	51.9	73.2	70.3	61.3	63.6	<b>91.1</b>
	Sad	65.0	54.9	53.5	70.7	69.2	61.1	62.4	88.4
	Angry	64.8	53.7	51.2	70.8	69.9	57.2	61.3	<b>88.6</b>
	Neutral	54.0	51.2	48.0	56.9	55.6	53.4	53.2	72.2
MCTN [37]	Happy	76.9	63.4	60.8	79.6	77.6	66.9	70.9	83.1
	Sad	76.7	64.4	60.4	78.9	77.1	68.6	71.0	82.8
	Angry	77.1	61.0	56.7	81.6	80.4	58.9	69.3	84.6
	Neutral	60.1	51.9	50.4	64.7	62.4	54.9	57.4	67.7
TransM [50]	Happy	78.4	64.5	61.1	81.6	80.2	66.5	72.1	85.5
	Sad	79.5	63.2	58.9	82.4	80.5	64.4	71.5	84.0
	Angry	81.0	65.0	60.7	83.9	81.7	66.9	73.2	86.1
	Neutral	60.2	49.9	50.7	65.2	62.4	52.4	56.8	67.1
SMIL [29]	Happy	80.5	66.5	63.8	83.1	81.8	68.2	74.0	86.8
	Sad	78.9	65.2	62.2	82.4	79.6	68.2	72.8	85.2
	Angry	79.6	67.2	61.8	83.1	82.0	67.8	73.6	84.9
	Neutral	60.2	50.4	48.8	65.4	62.2	52.6	56.6	68.9
GCNet [26]	Happy	81.9	67.3	66.6	83.7	82.5	69.8	75.3	87.7
	Sad	80.5	69.4	66.1	83.8	81.9	70.4	75.4	86.9
	Angry	80.1	66.2	64.2	82.5	81.6	68.1	73.8	85.2
	Neutral	61.8	51.1	49.6	66.2	63.5	53.3	57.6	71.1
CorrKD [24]	Happy	82.6	69.6	68.0	84.1	82.0	70.0	76.1	87.5
	Sad	82.7	<b>71.3</b>	67.6	83.4	82.2	72.5	76.6	85.9
	Angry	82.2	67.0	65.8	83.9	82.8	67.3	74.8	86.1
	Neutral	63.1	54.2	52.3	68.5	64.3	<b>57.2</b>	59.9	71.5
<b>HRLF (Ours)</b>	Happy	<b>84.9</b>	<b>71.8</b>	<b>69.7</b>	<b>86.4</b>	<b>85.6</b>	<b>72.3</b>	<b>78.5*</b>	88.1
	Sad	<b>83.7</b>	71.1	<b>69.0</b>	<b>85.3</b>	<b>83.9</b>	<b>73.6</b>	<b>77.8*</b>	86.4
	Angry	<b>83.4</b>	<b>69.1</b>	<b>67.2</b>	<b>84.5</b>	<b>83.5</b>	<b>70.9</b>	<b>76.4*</b>	86.7
	Neutral	<b>66.8</b>	<b>56.1</b>	<b>54.5</b>	<b>68.9</b>	<b>67.0</b>	56.9	<b>61.7*</b>	71.3

is eliminated, the worse performance demonstrates that aligning the high-level semantics in the representation by maximizing mutual information can generate favorable joint representations for the student network. (iii) Finally, we remove HAL, and the declined results illustrate that multi-scale adversarial learning can effectively align the representation distributions of student and teacher networks, thus effectively constraining the consistency across representations. This paradigm facilitates the recovery of missing semantics.

#### 4.5 Qualitative Analysis

To intuitively show the robustness of the proposed framework against modality missingness, we randomly select 100 samples in each emotion category on the IEMOCAP testing set to perform the visualization evaluation. The comparison models include CubeMLP [42] (complete-modality method), TransM [50] (joint learning-based missing-modality method), and GCNet [26] (generation-based missing-modality method). (i) As shown in Figure 6, CubeMLP fails to cope with the missing modality challenge because representations with different emotion categories are heavily confounded,

Table 3: Ablation results of inter-modality missingness case on the MOSI dataset.

Models	Testing Conditions							
	{ <i>l</i> }	{ <i>a</i> }	{ <i>v</i> }	{ <i>l, a</i> }	{ <i>l, v</i> }	{ <i>a, v</i> }	Avg.	{ <i>l, a, v</i> }
<b>HRLF (Full)</b>	<b>83.36</b>	<b>69.47</b>	<b>64.59</b>	<b>83.82</b>	<b>83.56</b>	<b>75.62</b>	<b>76.74</b>	<b>84.15</b>
w/o FRF	80.85	67.06	61.78	81.94	81.38	73.58	74.43	82.76
w/o HMI	81.54	67.72	62.70	82.45	81.90	74.22	75.09	83.25
w/o HAL	82.03	68.09	63.11	83.12	82.67	74.59	75.60	83.67

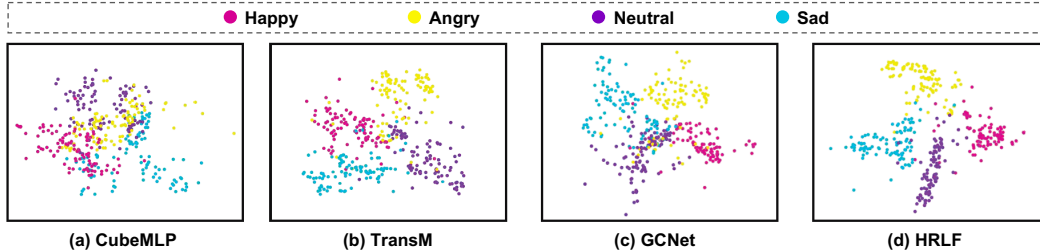


Figure 6: Visualization of representations from different methods with four emotion categories on the IEMOCAP testing set. The default testing conditions contain intra-modality missingness (*i.e.*, missing rate  $p = 0.5$ ) and inter-modality missingness (*i.e.*, only the language modality is available).

leading to the worst results. **(ii)** Although TransM and GCNet mitigate the indistinguishable emotion semantics to some extent, their performance is sub-optimal since the distribution boundaries of the different emotion representations are generally ambiguous and coupled. **(iii)** In comparison, our HRLF enables representations belonging to the same emotion category to form compact clusters, while representations of different categories are well separated. The above phenomenon benefits from the effective extraction of sentiment semantics and the precise filtering of task redundant information by the proposed hierarchical representation learning framework, which results in better joint multimodal representations. This further confirms the robustness and superiority of our framework.

## 5 Conclusion and Discussion

In this paper, we present a Hierarchical Representation Learning Framework (HRLF) to address diverse missing modality dilemmas in the MSA task. Specifically, we mine sentiment-relevant representations through a fine-grained representation factorization module. Additionally, the hierarchical mutual information maximization mechanism and the hierarchical adversarial learning mechanism are proposed for semantic and distributional alignment of representations of student and teacher networks to accurately reconstruct missing semantics and produce robust joint multimodal representations. Comprehensive experiments validate the superiority of our framework.

**Discussion of Limitation and Future Work.** The current method defines the modality missing cases as both inter-modality missingness and intra-modality missingness. Nevertheless, in real-world applications, modality missing cases may be very intricate and difficult to simulate. Consequently, the proposed method may suffer some minor performance loss when applied to real-world scenarios. In the future, we will explore more intricate modality missing cases and design suitable algorithms to compensate for this deficiency.

**Discussion of Broad Impacts.** The positive impact of our approach lies in the ability to significantly improve the robustness and stability of multimodal sentiment analysis systems against heterogeneous modality missingness in real-world applications. Nevertheless, this technology may have a negative impact when it falls into the wrong hands, *e.g.*, the proposed model is used for malicious purposes by injecting biased priors to recognize the emotions of specific groups.

## 6 Acknowledgements

This work was supported in part by National Key R&D Program of China 2021ZD0113502 and in part by Shanghai Municipal Science and Technology Major Project 2021SHZDZX0103.

## References

- [1] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9163–9171, 2019. 5
- [2] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. Openface: an open source facial behavior analysis toolkit. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–10. IEEE, 2016. 7
- [3] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. *Conference on Neural Information Processing Systems (NeurIPS)*, 29, 2016. 3
- [4] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42:335–359, 2008. 6
- [5] Jae Won Cho, Dong-Jin Kim, Jinsoo Choi, Yunjae Jung, and In So Kweon. Dealing with missing modalities in the visual question answer-difference prediction task through knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1592–1601, 2021. 3
- [6] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4794–4802, 2019. 3
- [7] Gilles Degottex, John Kane, Thomas Drugman, Tuomo Raitio, and Stefan Scherer. Covarep—a collaborative voice analysis repository for speech technologies. In *2014 Ieee International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 960–964. IEEE, 2014. 7
- [8] Changde Du, Changying Du, Hao Wang, Jinpeng Li, Wei-Long Zheng, Bao-Liang Lu, and Huiguang He. Semi-supervised deep generative modelling of incomplete multi-modality emotional data. In *Proceedings of the 26th ACM international conference on Multimedia (ACM MM)*, pages 108–116, 2018. 2
- [9] Yangtao Du, Dingkan Yang, Peng Zhai, Mingchen Li, and Lihua Zhang. Learning associative representation for facial expression recognition. In *IEEE International Conference on Image Processing (ICIP)*, pages 889–893, 2021. 1
- [10] Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. In *International Conference on Machine Learning (ICML)*, pages 1607–1616. PMLR, 2018. 3
- [11] Wei Han, Hui Chen, and Soujanya Poria. Improving multimodal fusion with hierarchical mutual information maximization for multimodal sentiment analysis. *arXiv preprint arXiv:2109.00412*, 2021. 2
- [12] Devamanyu Hazarika, Roger Zimmermann, and Soujanya Poria. Misa: Modality-invariant and-specific representations for multimodal sentiment analysis. In *Proceedings of the 28th ACM International Conference on Multimedia (ACM MM)*, pages 1122–1131, 2020. 1, 2, 4
- [13] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1921–1930, 2019. 3
- [14] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 3779–3787, 2019. 3
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 3, 5

- [16] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018. 5
- [17] Minhao Hu, Matthis Maillard, Ya Zhang, Tommaso Ciceri, Giammarco La Barbera, Isabelle Bloch, and Pietro Gori. Knowledge distillation from multi-modal to mono-modal segmentation networks. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23*, pages 772–781. Springer, 2020. 3, 5, 6
- [18] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning (ICML)*, pages 2649–2658. PMLR, 2018. 3
- [19] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 2018. 3
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7
- [21] Saurabh Kumar, Biplab Banerjee, and Subhasis Chaudhuri. Online sensor hallucination via knowledge distillation for multimodal image classification. *arXiv preprint arXiv:1908.10559*, 2019. 3, 5, 6
- [22] Mingcheng Li, Ding kang Yang, Yuxuan Lei, Shunli Wang, Shuaibing Wang, Liuzhen Su, Kun Yang, Yuzheng Wang, Mingyang Sun, and Lihua Zhang. A unified self-distillation framework for multimodal sentiment analysis with uncertain missing modalities. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 38, pages 10074–10082, 2024. 2
- [23] Mingcheng Li, Ding kang Yang, and Lihua Zhang. Towards robust multimodal sentiment analysis under uncertain signal missing. *IEEE Signal Processing Letters*, 30:1497–1501, 2023. 2
- [24] Mingcheng Li, Ding kang Yang, Xiao Zhao, Shuaibing Wang, Yan Wang, Kun Yang, Mingyang Sun, Dongliang Kou, Ziyun Qian, and Lihua Zhang. Correlation-decoupled knowledge distillation for multimodal sentiment analysis with incomplete modalities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12458–12468, 2024. 7, 8, 9
- [25] Yong Li, Yuanzhi Wang, and Zhen Cui. Decoupled multimodal distilling for emotion recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6631–6640, 2023. 1, 2, 3, 7, 8, 9
- [26] Zheng Lian, Lan Chen, Licai Sun, Bin Liu, and Jianhua Tao. Gcnet: graph completion network for incomplete multimodal learning in conversation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2, 7, 8, 9
- [27] Zhizhong Liu, Bin Zhou, Dianhui Chu, Yuhang Sun, and Lingqiang Meng. Modality translation-based multimodal sentiment analysis under uncertain missing modalities. *Information Fusion*, 101:101973, 2024. 2
- [28] Wei Luo, Mengying Xu, and Hanjiang Lai. Multimodal reconstruct and align net for missing modality problem in sentiment analysis. In *International Conference on Multimedia Modeling*, pages 411–422. Springer, 2023. 2
- [29] Mengmeng Ma, Jian Ren, Long Zhao, Sergey Tulyakov, Cathy Wu, and Xi Peng. Smil: Multimodal learning with severely missing modality. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pages 2302–2310, 2021. 2, 7, 8, 9
- [30] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 5191–5198, 2020. 3

- [31] Louis-Philippe Morency, Rada Mihalcea, and Payal Doshi. Towards multimodal sentiment analysis: Harvesting opinions from the web. In *Proceedings of the 13th International Conference on Multimodal Interfaces*, pages 169–176, 2011. 1
- [32] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International Conference on Machine Learning (ICML)*, pages 2642–2651. PMLR, 2017. 3
- [33] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3967–3976, 2019. 3
- [34] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 7
- [35] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5007–5016, 2019. 3
- [36] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. 6
- [37] Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-Philippe Morency, and Barnabás Póczos. Found in translation: Learning robust joint representations by cyclic translations between modalities. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 6892–6899, 2019. 2, 7, 8, 9
- [38] Masoomah Rahimpour, Jeroen Bertels, Ahmed Radwan, Henri Vandermeulen, Stefan Sunaert, Dirk Vandermeulen, Frederik Maes, Karolien Goffin, and Michel Koole. Cross-modal distillation to improve mri-based brain tumor segmentation with missing mri sequences. *IEEE Transactions on Biomedical Engineering*, 69(7):2153–2164, 2021. 3, 5, 6
- [39] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. 3
- [40] Roei Shraga, Haggai Roitman, Guy Feigenblat, and Mustafa Cannim. Web table retrieval using multimodal deep learning. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1399–1408, 2020. 1
- [41] Matthias Springstein, Eric Müller-Budack, and Ralph Ewerth. Quti! quantifying text-image consistency in multimodal documents. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2575–2579, 2021. 1
- [42] Hao Sun, Hongyi Wang, Jiaqing Liu, Yen-Wei Chen, and Lanfen Lin. Cubemlp: An mlp-based model for multimodal sentiment analysis and depression estimation. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 3722–3729, 2022. 2, 7, 8, 9
- [43] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019. 3
- [44] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2019, page 6558. NIH Public Access, 2019. 2
- [45] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1365–1374, 2019. 3
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 4

- [47] Hu Wang, Congbo Ma, Jianpeng Zhang, Yuan Zhang, Jodie Avery, Louise Hull, and Gustavo Carneiro. Learnable cross-modal knowledge distillation for multi-modal learning with missing modality. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 216–226. Springer, 2023. 3
- [48] Yansen Wang, Ying Shen, Zhun Liu, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Words can shift: Dynamically adjusting word representations using nonverbal behaviors. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 7216–7223, 2019. 6
- [49] Yuanzhi Wang, Zhen Cui, and Yong Li. Distribution-consistent modal recovering for incomplete multimodal learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22025–22034, 2023. 2
- [50] Zilong Wang, Zhaohong Wan, and Xiaojun Wan. Transmodality: An end2end fusion method with transformer for multimodal sentiment analysis. In *Proceedings of The Web Conference 2020*, pages 2514–2520, 2020. 2, 7, 8, 9
- [51] Wenke Xia, Xingjian Li, Andong Deng, Haoyi Xiong, Dejing Dou, and Di Hu. Robust cross-modal knowledge distillation for unconstrained videos. *arXiv preprint arXiv:2304.07775*, 2023. 3
- [52] Chenglin Yang, Lingxi Xie, Chi Su, and Alan L Yuille. Snapshot distillation: Teacher-student optimization in one generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2859–2868, 2019. 3
- [53] Dingkan Yang, Zhaoyu Chen, Yuzheng Wang, Shunli Wang, Mingcheng Li, Siao Liu, Xiao Zhao, Shuai Huang, Zhiyan Dong, Peng Zhai, and Lihua Zhang. Context de-confounded emotion recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19005–19015, June 2023. 1
- [54] Dingkan Yang, Shuai Huang, Haopeng Kuang, Yangtao Du, and Lihua Zhang. Disentangled representation learning for multimodal emotion recognition. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 1642–1651, 2022. 1, 3, 4
- [55] Dingkan Yang, Shuai Huang, Yang Liu, and Lihua Zhang. Contextual and cross-modal interaction for multi-modal speech emotion recognition. *IEEE Signal Processing Letters*, 29:2093–2097, 2022. 1
- [56] Dingkan Yang, Shuai Huang, Shunli Wang, Yang Liu, Peng Zhai, Liuzhen Su, Mingcheng Li, and Lihua Zhang. Emotion recognition for multiple context awareness. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*, pages 144–162. Springer, 2022. 1
- [57] Dingkan Yang, Haopeng Kuang, Shuai Huang, and Lihua Zhang. Learning modality-specific and-agnostic representations for asynchronous multimodal language sequences. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 1708–1717, 2022. 1, 3
- [58] Dingkan Yang, Haopeng Kuang, Kun Yang, Mingcheng Li, and Lihua Zhang. Towards asynchronous multimodal signal interaction and fusion via tailored transformers. *IEEE Signal Processing Letters*, 2024. 1
- [59] Dingkan Yang, Mingcheng Li, Linhao Qu, Kun Yang, Peng Zhai, Song Wang, and Lihua Zhang. Asynchronous multimodal video sequence fusion via learning modality-exclusive and-agnostic representations. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024. 1
- [60] Dingkan Yang, Mingcheng Li, Dongling Xiao, Yang Liu, Kun Yang, Zhaoyu Chen, Yuzheng Wang, Peng Zhai, Ke Li, and Lihua Zhang. Towards multimodal sentiment analysis debiasing via bias purification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 1

- [61] Ding kang Yang, Yang Liu, Can Huang, Mingcheng Li, Xiao Zhao, Yuzheng Wang, Kun Yang, Yan Wang, Peng Zhai, and Lihua Zhang. Target and source modality co-reinforcement for emotion understanding from asynchronous multimodal sequences. *Knowledge-Based Systems*, 265:110370, 2023. [1](#)
- [62] Ding kang Yang, Dongling Xiao, Ke Li, Yuzheng Wang, Zhaoyu Chen, Jinjie Wei, and Lihua Zhang. Towards multimodal human intention understanding debiasing via subject-deconfounding. *arXiv preprint arXiv:2403.05025*, 2024. [1](#)
- [63] Ding kang Yang, Kun Yang, Haopeng Kuang, Zhaoyu Chen, Yuzheng Wang, and Lihua Zhang. Towards context-aware emotion recognition debiasing from a causal demystification perspective via de-confounded training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. [1](#)
- [64] Ding kang Yang, Kun Yang, Mingcheng Li, Shunli Wang, Shuaibing Wang, and Lihua Zhang. Robust emotion recognition in context debiasing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12447–12457, 2024. [1](#)
- [65] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4133–4141, 2017. [3](#)
- [66] Wenmeng Yu, Hua Xu, Ziqi Yuan, and Jiele Wu. Learning modality-specific representations with self-supervised multi-task learning for multimodal sentiment analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pages 10790–10797, 2021. [1](#), [7](#), [8](#), [9](#)
- [67] Ziqi Yuan, Wei Li, Hua Xu, and Wenmeng Yu. Transformer-based feature reconstruction network for robust multimodal sentiment analysis. In *Proceedings of the 29th ACM International Conference on Multimedia (ACM MM)*, pages 4400–4407, 2021. [2](#)
- [68] Ziqi Yuan, Yihe Liu, Hua Xu, and Kai Gao. Noise imitation based adversarial training for robust multimodal sentiment analysis. *IEEE Transactions on Multimedia*, 2023. [2](#)
- [69] Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Tensor fusion network for multimodal sentiment analysis. *arXiv preprint arXiv:1707.07250*, 2017. [2](#)
- [70] Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Memory fusion network for multi-view sequential learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 32, 2018. [2](#)
- [71] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. Mosi: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos. *arXiv preprint arXiv:1606.06259*, 2016. [6](#), [7](#)
- [72] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, 2018. [6](#), [7](#)
- [73] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016. [3](#)
- [74] Jiandian Zeng, Tianyi Liu, and Jiantao Zhou. Tag-assisted multimodal sentiment analysis under uncertain missing modalities. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1545–1554, 2022. [2](#)
- [75] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2881–2890, 2017. [3](#)
- [76] Jinming Zhao, Ruichen Li, and Qin Jin. Missing modality imagination network for emotion recognition with uncertain missing modalities. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2608–2618, 2021. [2](#)

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Please refer to the "Abstract" and "1 Introduction" for our paper's contributions and scopes.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Please refer to the "5 Conclusion and Discussion" section for the limitations of our work.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?



Answer: [NA]

Justification: No theory assumptions and proofs are provided in the paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The “4.2 Implementation Details” section of the paper describes all the information needed to reproduce the main experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: The paper does not provide open access to the data and code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The “4.2 Implementation Details” section of the paper specify all the training and testing details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: In Tables 1 and 2 of the paper, we conducted significance tests on the experimental results to demonstrate the superior performance of the proposed framework.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The “4.2 Implementation Details” section of the paper explains that all experiments are conducted on four NVIDIA Tesla V100 GPUs.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn’t make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Our research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Please refer to the “5 Conclusion and Discussion” sections for the broader impacts of our work

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The MOSI, MOSEI and IEMOCAP datasets and the Pytorch toolbox in this paper are existing assets and we cite the references.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.