

Correlation-Decoupled Knowledge Distillation for Multimodal Sentiment Analysis with Incomplete Modalities

Mingcheng Li^{1,2}[§] Dingkang Yang^{1,2}[§] Xiao Zhao^{1,2} Shuaibing Wang^{1,2} Yan Wang¹
 Kun Yang¹ Mingyang Sun^{1,2} Dongliang Kou^{1,2} Ziyun Qian^{1,2} Lihua Zhang^{1,2,3,4}[§]

¹Academy for Engineering and Technology, Fudan University ²Cognition and Intelligent Technology Laboratory (CIT Lab)

³Jilin Provincial Key Laboratory of Intelligence Science and Engineering, Changchun, China

⁴Engineering 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) aims to understand human sentiment through multimodal data. Most MSA efforts are based on the assumption of modality completeness. However, in real-world applications, some practical factors cause uncertain modality missingness, which drastically degrades the model’s performance. To this end, we propose a Correlation-decoupled Knowledge Distillation (CorrKD) framework for the MSA task under uncertain missing modalities. Specifically, we present a sample-level contrastive distillation mechanism that transfers comprehensive knowledge containing cross-sample correlations to reconstruct missing semantics. Moreover, a category-guided prototype distillation mechanism is introduced to capture cross-category correlations using category prototypes to align feature distributions and generate favorable joint representations. Eventually, we design a response-disentangled consistency distillation strategy to optimize the sentiment decision boundaries of the student network through response disentanglement and mutual information maximization. Comprehensive experiments on three datasets indicate that our framework can achieve favorable improvements compared with several baselines.

1. Introduction

“Correlations serve as the beacon through the fog of the missingness.”

–Lee & Dicken

Multimodal sentiment analysis (MSA) has attracted wide attention in recent years. Different from the traditional unimodal-based emotion recognition task [7], MSA

[§]Corresponding author. [§]Equal contribution.


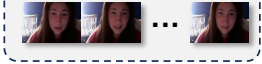

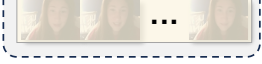
Modality	Content	Prediction	Label
Language	It was a great movie and I loved it.	Positive ✓	Positive
Audio			
Visual			
Language	It was a great movie and I loved it.	Neutral ✗	Positive
Audio			
Visual			

Figure 1. Traditional model outputs correct prediction when inputting the sample with complete modalities, but incorrectly predicts the sample with missing modalities. We define two missing modality cases: (i) intra-modality missingness (*i.e.*, the pink areas) and (ii) inter-modality missingness (*i.e.*, the yellow area).

understands and recognizes human emotions through multiple modalities, including language, audio, and visual [28]. Previous studies have shown that combining complementary information among different modalities facilitates the generation of more valuable joint multimodal representations [34, 36]. Under the deep learning paradigm [3, 17, 42, 43, 54, 59, 60], numerous studies assuming the availability of all modalities during both training and inference stages [10, 19, 22, 49–53, 55–58, 62]. Nevertheless, this assumption often fails to align with real-world scenarios, where factors such as background noise, sensor constraints, and privacy concerns may lead to uncertain modality missingness issues. Modality missingness can significantly impair the effectiveness of well-trained models based on complete modalities. For instance, as shown in Figure 1, the entire visual modality is missing, and some frame-level fea-

tures in the language and audio modalities are missing, leading to an incorrect sentiment prediction.

In recent years, many works [20, 21, 23, 24, 32, 45, 46, 66] attempt to address the problem of missing modalities in MSA. As a typical example, MCTN [32] guarantees the model’s robustness to the missing modality case by learning a joint representation through cyclic translation from the source modality to the target modality. However, these methods suffer from the following limitations: (i) inadequate interactions based on individual samples lack the mining of holistically structured semantics. (ii) Failure to model cross-category correlations leads to loss of sentiment-relevant information and confusing distributions among categories. (iii) Coarse supervision ignores the semantic and distributional alignment.

To address the above issues, we present a **Correlation-decoupled Knowledge Distillation (CorrKD)** framework for the MSA task under uncertain missing modalities. There are three core contributions in CorrKD based on the tailored components. Specifically, (i) the proposed sample-level contrastive distillation mechanism captures the holistic cross-sample correlations and transfers valuable supervision signals via sample-level contrastive learning. (ii) Meanwhile, we design a category-guided prototype distillation mechanism that leverages category prototypes to transfer intra- and inter-category feature variations, thus delivering sentiment-relevant information and learning robust joint multimodal representations. (iii) Furthermore, we introduce a response-disentangled consistency distillation strategy to optimize sentiment decision boundaries and encourage distribution alignment by decoupling heterogeneous responses and maximizing mutual information between homogeneous sub-responses. Based on these components, CorrKD significantly improves MSA performance under uncertain missing-modality and complete-modality testing conditions on three multimodal benchmarks.

2. Related Work

2.1. Multimodal Sentiment Analysis

MSA aims to understand and analyze human sentiment utilizing multiple modalities. Mainstream MSA studies [9, 10, 22, 37, 50, 53, 55–58] focus on designing complex fusion paradigms and interaction mechanisms to enhance the performance of sentiment recognition. For instance, CubeMLP [37] utilizes three independent multi-layer perceptron units for feature-mixing on three axes. However, these approaches based on complete modalities cannot be deployed in real-world applications. Mainstream solutions for the missing modality problem can be summarized in two categories: (i) generative methods [6, 23, 25, 45] and (ii) joint learning methods [24, 32, 46, 66]. Reconstruction methods generate missing features and seman-

tics in modalities based on available modalities. For example, TFR-Net [63] leverages the feature reconstruction module to guide the extractor to reconstruct missing semantics. MVAE [6] solves the modality missing problem by the semi-supervised multi-view deep generative framework. Joint learning efforts refer to learning joint multimodal representations utilizing correlations among modalities. For instance, MMIN [69] generates robust joint multimodal representations via cross-modality imagination. TATE [66] presents a tag encoding module to guide the network to focus on missing modalities. However, the aforementioned approaches fail to account for the correlations among samples and categories, leading to inadequate compensation for the missing semantics in modalities. In contrast, we design effective learning paradigms to adequately capture potential inter-sample and inter-category correlations.

2.2. Knowledge Distillation

Knowledge distillation utilizes additional supervisory information from the pre-trained teacher’s network to assist in the training of the student’s network [11]. Knowledge distillation methods can be roughly categorized into two types, distillation from intermediate features [15, 29, 38, 61] and responses [4, 8, 27, 48, 68]. Many studies [13, 18, 33, 40, 47] employ knowledge distillation for MSA tasks with missing modalities. The core concept of these efforts is to transfer “dark knowledge” from teacher networks trained by complete modalities to student networks trained by missing modalities. The teacher model typically produces more valuable feature presentations than the student model. For instance, [13] utilizes the complete-modality teacher network to implement supervision on the unimodal student network at both feature and response levels. Despite promising outcomes, they are subject to several significant limitations: (i) Knowledge transfer is limited to individual samples, overlooking the exploitation of clear correlations among samples and among categories. (ii) Supervision on student networks is coarse-grained and inadequate, without considering the potential alignment of feature distributions. To this end, we propose a correlation-decoupled knowledge distillation framework that facilitates the learning of robust joint representations by refining and transferring the cross-sample, cross-category, and cross-target correlations.

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. $T_m(\cdot)$ is the sequence length and $d_m(\cdot)$ is the embedding dimension, where $m \in \{L, A, V\}$. Meanwhile, the incomplete modality is denoted

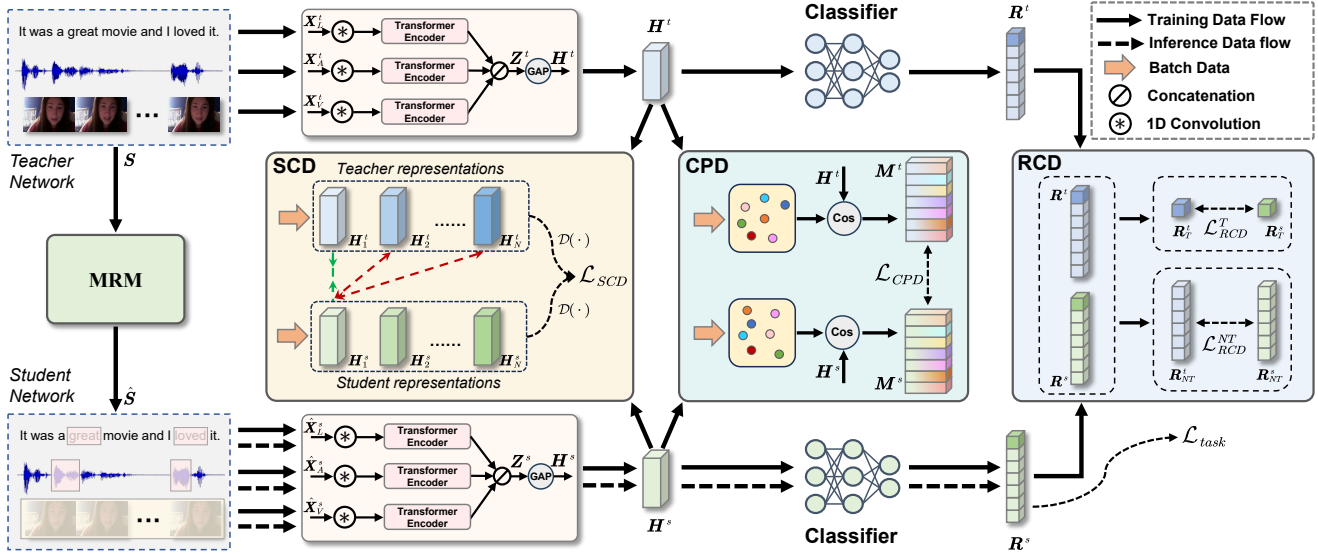


Figure 2. The structure of our CorrKD, which consists of three core components: Sample-level Contrastive Distillation (SCD) mechanism, Category-guided Prototype Distillation (CPD) mechanism, and Response-disentangled Consistency Distillation (RCD) strategy.

as \hat{X}_m . We define two missing modality cases to simulate the most natural and holistic challenges in real-world scenarios: (i) *intra-modality missingness*, which indicates some frame-level features in the modality sequences are missing. (ii) *inter-modality missingness*, which denotes some modalities are entirely missing. Our goal is to recognize the utterance-level sentiments by utilizing the multimodal data with missing modalities.

3.2. Overall Framework

Figure 2 illustrates the main workflow of CorrKD. The teacher network and the student network adopt a consistent structure but have different parameters. During the training phase, our CorrKD procedure is as follows: (i) we train the teacher network with complete-modality samples and then freeze its parameters. (ii) Given a video segment sample S , we generate a missing-modality sample \hat{S} with the Modality Random Missing (MRM) strategy. MRM simultaneously performs intra-modality missing and inter-modality missing, and the raw features of the missing portions are replaced with zero vectors. S and \hat{S} are fed into the initialized student network and the trained teacher network, respectively. (iii) We input the samples S and \hat{S} into the modality representation fusion module to obtain the joint multimodal representations H^t and H^s . (iv) The sample-level contrastive distillation mechanism and the category-guided prototype distillation mechanism are utilized to learn the feature consistency of H^t and H^s . (v) These representations are fed into the task-specific fully-connected layers and the softmax function to obtain the network responses R^t and R^s . (vi) The response-disentangled consistency distillation

strategy is applied to maintain consistency in the response distribution, and then R^s is used to perform classification. In the inference phase, testing samples are only fed into the student network for downstream tasks. Subsequent sections provide details of the proposed components.

3.3. Modality Representation Fusion

We introduce the extraction and fusion processes of modality representations using the student network as an example. The incomplete modality $\hat{X}_m^s \in \mathbb{R}^{T_m \times d_m}$ with $m \in \{L, A, V\}$ is fed into the student network. Firstly, \hat{X}_m^s passes through a 1D temporal convolutional layer with kernel size 3×3 and adds the positional embedding [39] to obtain the preliminary representations, denoted as $\hat{F}_m^s = \mathbf{W}_{3 \times 3}(\hat{X}_m^s) + PE(T_m, d) \in \mathbb{R}^{T_m \times d}$. Each F_m^s is fed into a Transformer [39] encoder $\mathcal{F}_\phi^s(\cdot)$, capturing the modality dynamics of each sequence through the self-attention mechanism to yield representations E_m^s , denoted as $E_m^s = \mathcal{F}_\phi^s(F_m^s)$. The representations E_m^s are concatenated to obtain Z^s , expressed as $Z^s = [E_L^s, E_A^s, E_V^s] \in \mathbb{R}^{T_m \times 3d}$. Subsequently, Z^s is fed into the Global Average Pooling (GAP) to further enhance and refine the features, yielding the joint multimodal representation $H^s \in \mathbb{R}^{3d}$. Similarly, the joint multimodal representation generated by the teacher network is represented as $H^t \in \mathbb{R}^{3d}$.

3.4. Sample-level Contrastive Distillation

Most previous studies of MSA tasks with missing modalities [33, 40, 47] are sub-optimal, exploiting only one-sided information within a single sample and neglecting to consider comprehensive knowledge across samples. To

this end, we propose a Sample-level Contrastive Distillation (SCD) mechanism that enriches holistic knowledge encoding by implementing contrastive learning between sample-level representations of student and teacher networks. This paradigm prompts models to sufficiently capture intra-sample dynamics and inter-sample correlations to generate and transfer valuable supervision signals, thus precisely recovering the missing semantics. The rationale of SCD is to take contrastive learning within all mini-batches, constraining the representations in two networks originating from the same sample to be similar, and the representations originating from different samples to be distinct.

Specifically, given a mini-batch with N samples $\mathbf{B} = \{\mathbf{S}_0, \mathbf{S}_1, \dots, \mathbf{S}_N\}$, we obtain their sets of joint multimodal representations in teacher and student networks, denoted as $\{\mathbf{H}_1^w, \mathbf{H}_2^w, \dots, \mathbf{H}_N^w\}$ with $w \in \{t, s\}$. For the same input sample, we narrow the distance between the joint representations of the teacher and student networks and enlarge the distance between the representations for different samples. The contrastive distillation loss is formulated as follows:

$$\mathcal{L}_{SCD} = \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathcal{D}(\mathbf{H}_i^s, \mathbf{H}_i^t)^2 + \max\{0, \eta - \mathcal{D}(\mathbf{H}_i^s, \mathbf{H}_j^t)\}^2, \quad (1)$$

where $\mathcal{D}(\mathbf{H}^s, \mathbf{H}^t) = \|\mathbf{H}^s - \mathbf{H}^t\|_2$, $\|\cdot\|_2$ represents ℓ_2 norm function, and η is the predefined distance boundary. When negative pairs are distant enough (*i.e.*, greater than boundary η), the loss is set to 0, allowing the model to focus on other pairs. Since the sample-level representation contains holistic emotion-related semantics, such a contrastive objective facilitates the student network to learn more valuable knowledge from the teacher network.

3.5. Category-guided Prototype Distillation

MSA data usually suffers from the dilemmas of high intra-category diversity and high inter-category similarity. Previous approaches [13, 18, 33] based on knowledge distillation to address the modality missing problem simply constrain the feature consistency of the teacher and student networks. The rough manner lacks consideration of cross-category correlation and feature variations, leading to ambiguous feature distributions. To this end, we propose a Category-guided Prototype Distillation (CPD) mechanism, with the core insight of refining and transferring knowledge of intra- and inter-category feature variations via category prototypes, which is widely utilized in the field of few-shot learning [35]. The category prototype represents the embedding center of every sentiment category, denoted as:

$$\mathbf{c}_k = \frac{1}{|\mathbf{B}_k|} \sum_{\mathbf{S}_i \in \mathbf{B}_k} \mathbf{H}_i, \quad (2)$$

where \mathbf{B}_k denotes the set of samples labeled with category k in the mini-batch, and \mathbf{S}_i denotes the i -th sample in \mathbf{B}_k .

The intra- and inter-category feature variation of the sample \mathbf{S}_i is defined as follows:

$$\mathbf{M}_k(i) = \frac{\mathbf{H}_i \mathbf{c}_k^\top}{\|\mathbf{H}_i\|_2 \|\mathbf{c}_k\|_2}, \quad (3)$$

where $\mathbf{M}_k(i)$ denotes the similarity between the sample \mathbf{S}_i and the prototype \mathbf{c}_k . If the sample \mathbf{S}_i is of category k , $\mathbf{M}_k(i)$ represents intra-category feature variation. Otherwise, it represents inter-category feature variation. The teacher and student networks compute similarity matrices \mathbf{M}^t and \mathbf{M}^s , respectively. We minimize the squared Euclidean distance between the two similarity matrices to maintain the consistency of two multimodal representations. The prototype distillation loss is formulated as:

$$\mathcal{L}_{CPD} = \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K \|\mathbf{M}_k^s(i) - \mathbf{M}_k^t(i)\|_2, \quad (4)$$

where K is the category number of the mini-batch.

3.6. Response-disentangled Consistency Distillation

Most knowledge distillation studies [15, 29, 38, 61] focus on extracting knowledge from intermediate features of networks. Although the model's response (*i.e.*, the predicted probability of the model's output) presents a higher level of semantics than the intermediate features, response-based methods achieve significantly worse performance than feature-based methods [41]. Inspired by [67], the model's response consists of two parts: (i) Target Category Response (TCR), which represents the prediction of the target category and describes the difficulty of identifying each training sample. (ii) Non-Target Category Response (NTCR), which denotes the prediction of the non-target category and reflects the decision boundaries of the remaining categories to some extent. The effects of TCR and NTCR in traditional knowledge distillation loss are coupled, *i.e.*, high-confidence TCR leads to low-impact NTCR, thus inhibiting effective knowledge transfer. Consequently, we disentangle the heterogeneous responses and constrain the consistency between the homogeneous responses. From the perspective of information theory, knowledge consistency between responses can be characterized as maintaining high mutual information between teacher and student networks [1]. This schema captures beneficial semantics and encourages distributional alignment.

Specifically, the joint multimodal representation \mathbf{H}^w with $w \in \{t, s\}$ of teacher and student networks pass through fully-connected layers and softmax function to obtain response \mathbf{R}^w . Based on the target indexes, we decouple the response \mathbf{R}^w to obtain TCR \mathbf{R}_T^w and NTCR \mathbf{R}_{NT}^w . Define $\mathbf{Q} \in \mathcal{Q}$ and $\mathbf{U} \in \mathcal{U}$ as two random variables. Formulaically, the marginal probability density functions of \mathbf{Q} and \mathbf{U} are denoted as $P(\mathbf{Q})$ and $P(\mathbf{U})$. $P(\mathbf{Q}, \mathbf{U})$ is regarded as the joint probability density function. The mutual

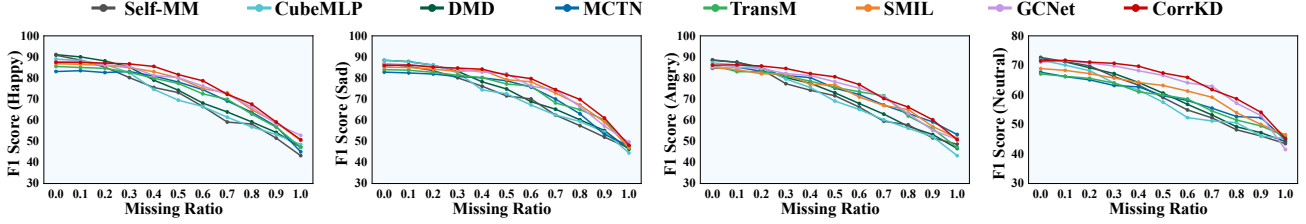


Figure 3. Comparison results of intra-modality missingness on IEMOCAP. We comprehensively report the F1 score for the happy, sad, angry, and neutral categories at various missing ratios.

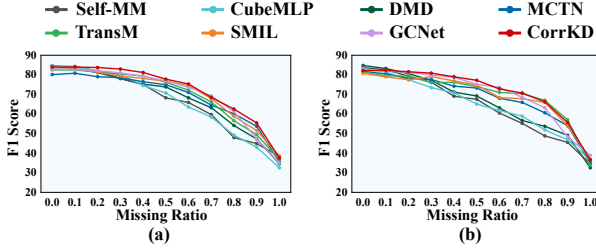


Figure 4. Comparison results of intra-modality missingness on (a) MOSI and (b) MOSEI. We report the F1 score at various ratios.

information between Q and U is represented as follows:

$$I(Q, U) = \int_{\mathcal{Q}} \int_{\mathcal{U}} P(Q, U) \log \left(\frac{P(Q, U)}{P(Q)P(U)} \right) dQ dU. \quad (5)$$

The mutual information $I(Q, U)$ can be written as the Kullback-Leibler divergence between the joint probability distribution \mathbb{P}_{QU} and the product of the marginal distributions $\mathbb{P}_Q \mathbb{P}_U$, denoted as $I(Q, U) = D_{KL}(\mathbb{P}_{QU} \| \mathbb{P}_Q \mathbb{P}_U)$. For efficient and stable computation, the Jensen-Shannon divergence [12] is employed in our case to estimate the mutual information, which is denoted as follows:

$$\begin{aligned} I(Q, U) &\geq \hat{I}_{\theta}^{(JSD)}(Q, U) \\ &= \mathbb{E}_{P(Q, U)} \left[-\log \left(1 + e^{-\mathcal{F}_{\theta}(Q, U)} \right) \right] \\ &\quad - \mathbb{E}_{P(Q)P(U)} \left[\log \left(1 + e^{\mathcal{F}_{\theta}(Q, U)} \right) \right], \end{aligned} \quad (6)$$

where $\mathcal{F}_{\theta} : Q \times U \rightarrow \mathbb{R}$ is formulated as an instantiated statistical network with parameters θ . We only need to maximize the mutual information without focusing on its precise value. Consequently, the distillation loss based on the mutual information estimation is formatted as follows:

$$\mathcal{L}_{RCD} = \mathcal{L}_{RCD}^T + \mathcal{L}_{RCD}^{NT} = -I(\mathbf{R}_T^t, \mathbf{R}_T^s) - I(\mathbf{R}_{NT}^t, \mathbf{R}_{NT}^s). \quad (7)$$

Finally, the overall training objective \mathcal{L}_{total} is expressed as $\mathcal{L}_{total} = \mathcal{L}_{task} + \mathcal{L}_{SCD} + \mathcal{L}_{CPD} + \mathcal{L}_{RCD}$, where \mathcal{L}_{task} is the standard cross-entropy loss.

4. Experiments

4.1. Datasets and Evaluation Metrics

We conduct extensive experiments on three MSA datasets with word-aligned data, including MOSI [64], MOSEI [65], and IEMOCAP [2]. **MOSI** is a realistic dataset that comprises 2,199 short monologue video clips. There are 1,284, 229, and 686 video clips in train, valid, and test data, respectively. **MOSEI** is a dataset consisting of 22,856 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 labeled by human annotators with a sentiment score of -3 (strongly negative) to +3 (strongly positive). On the MOSI and MOSEI datasets, we utilize weighted F1 score computed for positive/negative classification results as evaluation metrics. **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 [44], four emotions (*i.e.*, happy, sad, angry, and neutral) are selected for emotion recognition. For evaluation, we report the F1 score for each category.

4.2. Implementation Details

Feature Extraction. The Glove embedding [31] 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 toolkit [5] to extract 74-dimensional acoustic features, including 12 Mel-frequency cepstral coefficients (MFCCs), voiced/unvoiced segmenting features, and glottal source parameters. For the visual modality, we utilize the Facet [14] to indicate 35 facial action units, recording facial movement to express emotions.

Experimental Setup. All models are built on the Pytorch [30] toolbox with NVIDIA Tesla V100 GPUs. The Adam optimizer [16] is employed for network optimization. For MOSI, MOSEI, and IEMOCAP, the detailed hyper-parameter settings are as follows: the learning rates are $\{4e-3, 2e-3, 4e-3\}$, the batch sizes are $\{64, 32, 64\}$, the epoch numbers are $\{50, 20, 30\}$, the attention heads are $\{10, 8, 10\}$, and the distance boundaries η are $\{1.2, 1.0, 1.4\}$. The embedding dimension is 40 on all three datasets. The hyper-parameters are determined via

Table 1. Comparison results under inter-modality missing and complete-modality testing conditions on MOSI and MOSEI.

Dataset	Models	Testing Conditions							
		$\{l\}$	$\{a\}$	$\{v\}$	$\{l, a\}$	$\{l, v\}$	$\{a, v\}$	Avg.	$\{l, a, v\}$
MOSI	Self-MM [62]	67.80	40.95	38.52	69.81	74.97	47.12	56.53	84.64
	CubeMLP [37]	64.15	38.91	43.24	63.76	65.12	47.92	53.85	84.57
	DMD [22]	68.97	43.33	42.26	70.51	68.45	50.47	57.33	84.50
	MCTN [32]	75.21	59.25	58.57	77.81	74.82	64.21	68.31	80.12
	TransM [46]	77.64	63.57	56.48	82.07	80.90	67.24	71.32	82.57
	SMIL [26]	78.26	67.69	59.67	79.82	79.15	71.24	72.64	82.85
	GCNet [23]	80.91	65.07	58.70	84.73	83.58	70.02	73.84	83.20
	CorrKD	81.20	66.52	60.72	83.56	82.41	73.74	74.69	83.94
MOSEI	Self-MM [62]	71.53	43.57	37.61	75.91	74.62	49.52	58.79	83.69
	CubeMLP [37]	67.52	39.54	32.58	71.69	70.06	48.54	54.99	83.17
	DMD [22]	70.26	46.18	39.84	74.78	72.45	52.70	59.37	84.78
	MCTN [32]	75.50	62.72	59.46	76.64	77.13	64.84	69.38	81.75
	TransM [46]	77.98	63.68	58.67	80.46	78.61	62.24	70.27	81.48
	SMIL [26]	76.57	65.96	60.57	77.68	76.24	66.87	70.65	80.74
	GCNet [23]	80.52	66.54	61.83	81.96	81.15	69.21	73.54	82.35
	CorrKD	80.76	66.09	62.30	81.74	81.28	71.92	74.02	82.16

the validation set. The raw features at the modality missing positions are replaced by zero vectors. To ensure an equitable comparison, we re-implement the state-of-the-art (SOTA) methods using the publicly available codebases 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 compare CorrKD with seven representative and reproducible SOTA methods, including complete-modality methods: Self-MM [62], CubeMLP [37], and DMD [22], and missing-modality methods: 1) joint learning methods (*i.e.*, MCTN [32] and TransM [46]), and 2) generative methods (*i.e.*, SMIL [26] and GCNet [23]). Extensive experiments are implemented to thoroughly evaluate the robustness and effectiveness of CorrKD in the cases of intra-modality and inter-modality missingness.

Robustness to Intra-modality Missingness. We randomly drop frame-level features in modality sequences with ratio $p \in \{0.1, 0.2, \dots, 1.0\}$ to simulate testing conditions of intra-modality missingness. Figures 3 and 4 show the performance curves of models with various p values, which intuitively reflect the model’s robustness. We have the following important observations. (i) As the ratio p increases, the performance of all models decreases. This phenomenon demonstrates that intra-modality missingness leads to a considerable loss of sentiment semantics and fragile joint multimodal representations. (ii) Compared to the complete-modality methods (*i.e.*, Self-MM, CubeMLP,

and DMD), our CorrKD achieves significant performance advantages in the missing-modality testing conditions and competitive performance in the complete-modality testing conditions. The reason is that complete-modality methods are based on the assumption of data completeness, whereas customized training paradigms for missing modalities perform better at capturing and reconstructing valuable sentiment semantics from incomplete multimodal data. (iii) Compared to the missing-modality methods, our CorrKD exhibits the strongest robustness. Benefiting from the decoupling and modeling of inter-sample, inter-category, and inter-response correlations by the proposed correlation decoupling schema, the student network acquires informative knowledge to reconstruct valuable missing semantics and produces robust multimodal representations.

Robustness to Inter-modality Missingness. In Table 1 and 2, we drop some entire modalities in the samples to simulate testing conditions of inter-modality missingness. The notation “ $\{l\}$ ” indicates that only the language modality is available, while audio and visual modalities are missing. “ $\{l, a, v\}$ ” represents the complete-modality testing condition where all modalities are available. “Avg.” indicates the average performance across six missing-modality testing conditions. We present the following significant insights. (i) Inter-modality missingness causes performance degradation for all models, suggesting that the integration of complementary information from heterogeneous modalities enhances the sentiment semantics within joint representations. (ii) In the testing conditions of the inter-modality missingness, our CorrKD has superior performance among

Table 2. Comparison results under six testing conditions of inter-modality missingness and the complete-modality condition on IEMOCAP.

Models	Categories	Testing Conditions							
		{l}	{a}	{v}	{l, a}	{l, v}	{a, v}	Avg.	{l, a, v}
Self-MM [62]	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	72.7
CubeMLP [37]	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	88.5
	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 [22]	Happy	69.5	55.4	51.9	73.2	70.3	61.3	63.6	91.1
	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	88.6
	Neutral	54.0	51.2	48.0	56.9	55.6	53.4	53.2	72.2
MCTN [32]	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 [46]	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 [26]	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 [23]	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	Happy	82.6	69.6	68.0	84.1	82.0	70.0	76.1	87.5
	Sad	82.7	71.3	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	57.2	59.9	71.5

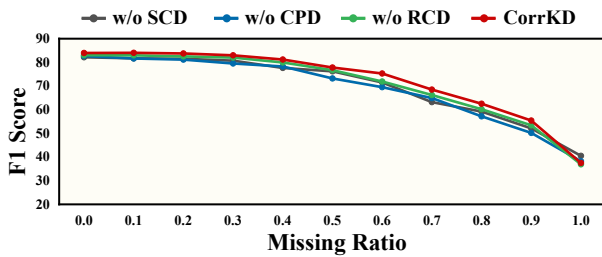


Figure 5. Ablation results of intra-modality missingness using various missing ratios on MOSI.

the majority of metrics, proving its strong robustness. For example, on the MOSI dataset, CorrKD’s average F1 score is improved by 0.85% compared to GCNet, and in particular by 3.72% in the testing condition where language modality is missing (*i.e.*, {a, v}). The merit stems from the pro-

Table 3. Ablation results for the testing conditions of inter-modality missingness on MOSI.

Models	Testing Conditions							
	{l}	{a}	{v}	{l, a}	{l, v}	{a, v}	Avg.	{l, a, v}
CorrKD	81.20	66.52	60.72	83.56	82.41	73.74	74.69	83.94
w/o SCD	78.80	64.96	57.49	81.95	80.53	71.05	72.46	82.13
w/o CPD	79.23	63.72	57.83	80.11	79.45	70.53	71.81	82.67
w/o RCD	79.73	65.32	59.21	82.14	81.05	72.18	73.27	83.05

posed framework’s capability of decoupling and modeling potential correlations at multiple levels to capture discriminative and holistic sentiment semantics. (iii) In the unimodal testing conditions, the performance of CorrKD with only the language modality favorably outperforms other cases, with comparable results to the complete-modality

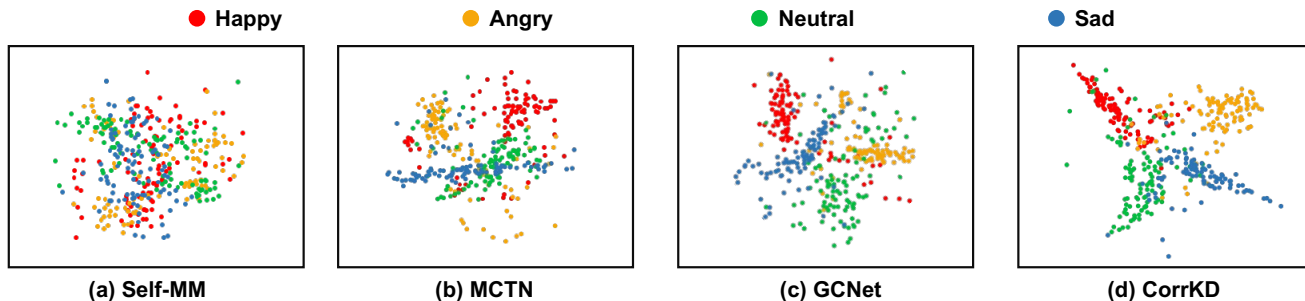


Figure 6. Visualization of representations from different methods with four emotion categories on the IEMOCAP testing set. The default testing conditions contain intra-modality missingness (*i.e.*, missing ratio $p = 0.5$) and inter-modality missingness (*i.e.*, only the language modality is available). The red, orange, green, and blue markers represent the happy, angry, neutral, and sad emotions, respectively.

case. In the bimodal testing conditions, cases containing the language modality perform the best, even surpassing the complete-modality case in individual metrics. This phenomenon proves that language modality encompasses the richest knowledge information and dominates the sentiment inference and missing semantic reconstruction.

4.4. Ablation Studies

To validate the effectiveness and necessity of the proposed mechanisms and strategies in CorrKD, we conduct ablation studies under two missing-modality cases on the MOSI dataset, as shown in Table 3 and Figure 5. The principal findings are outlined as follows. (i) When SCD is eliminated, there is a noticeable degradation in model performance under both missing cases. This phenomenon suggests that mining and transferring comprehensive cross-sample correlations is essential for recovering missing semantics in student networks. (ii) The worse results under the two missing modality scenarios without CPD indicate that capturing cross-category feature variations and correlations facilitates deep alignment of feature distributions between both networks to produce robust joint multimodal representations. (iii) Moreover, we substitute the KL divergence loss for the proposed RCD. The declining performance gains imply that decoupling heterogeneous responses and maximizing mutual information between homogeneous responses motivate the student network to adequately reconstruct meaningful sentiment semantics.

4.5. Qualitative Analysis

To intuitively show the robustness of the proposed framework against modality missingness, we randomly choose 100 samples from each emotion category on the IEMOCAP testing set for visualization analysis. The comparison models include Self-MM [62] (*i.e.*, complete-modality method), MCTN [32] (*i.e.*, joint learning-based missing-modality method), and GCNet [23] (*i.e.*, generative-based missing-modality method). (i) As shown in Figure 6, Self-

MM cannot address the modality missing challenge, as the representations of different emotion categories are heavily confounded, leading to the least favorable outcomes. (ii) Although MCTN and GCNet somewhat alleviate the issue of indistinct emotion semantics, their effectiveness remains limited since the distribution boundaries of the different emotion representations are generally ambiguous and coupled. (iii) Conversely, our CorrKD ensures that representations of the same emotion category form compact clusters, while representations of different categories are clearly separated. These observations confirm the robustness and superiority of our framework, as it sufficiently decouples inter-sample, inter-category and inter-response correlations.

5. Conclusions

In this paper, we present a correlation-decoupled knowledge distillation framework (CorrKD) to address diverse missing modality dilemmas in the MSA task. Concretely, we propose a sample-level contrast distillation mechanism that utilizes contrastive learning to capture and transfer cross-sample correlations to precisely reconstruct missing semantics. Additionally, we present a category-guided prototype distillation mechanism that learns cross-category correlations through category prototypes, refining sentiment-relevant semantics for improved joint representations. Eventually, a response-disentangled consistency distillation is proposed to encourage distribution alignment between teacher and student networks. Extensive experiments confirm the effectiveness of our framework.

Acknowledgements

This work is supported in part by the Shanghai Municipal Science and Technology Committee of Shanghai Outstanding Academic Leaders Plan (No. 21XD1430300), and in part by the National Key R&D Program of China (No. 2021ZD0113503).

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. 4
- [2] 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. 5
- [3] Jiawei Chen, Dingkang Yang, Yue Jiang, Yuxuan Lei, and Lihua Zhang. Miss: A generative pretraining and finetuning approach for med-vqa. *arXiv preprint arXiv:2401.05163*, 2024. 1
- [4] 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. 2
- [5] Gilles Degottex, John Kane, Thomas Drugman, Tuomo Raitio, and Stefan Scherer. Covarep—a collaborative voice analysis repository for speech technologies. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 960–964. IEEE, 2014. 5
- [6] Changde Du, Changying Du, Hao Wang, Jinpeng Li, Weilong Zheng, Bao-Liang Lu, and Huiguang He. Semi-supervised deep generative modelling of incomplete multimodality emotional data. In *Proceedings of the 26th ACM international conference on Multimedia (ACM MM)*, pages 108–116, 2018. 2
- [7] Yangtao Du, Dingkang 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
- [8] 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. 2
- [9] 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
- [10] 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
- [11] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2
- [12] 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
- [13] 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. 2, 4
- [14] iMotions. Facial expression analysis. 2017. 5
- [15] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 2018. 2, 4
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [17] Haopeng Kuang, Dingkang Yang, Shunli Wang, Xiaoying Wang, and Lihua Zhang. Towards simultaneous segmentation of liver tumors and intrahepatic vessels via cross-attention mechanism. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023. 1
- [18] Saurabh Kumar, Biplab Banerjee, and Subhasis Chaudhuri. Online sensor hallucination via knowledge distillation for multimodal image classification. *arXiv preprint arXiv:1908.10559*, 2019. 2, 4
- [19] Yuxuan Lei, Dingkang Yang, Mingcheng Li, Shunli Wang, Jiawei Chen, and Lihua Zhang. Text-oriented modality reinforcement network for multimodal sentiment analysis from unaligned multimodal sequences. *arXiv preprint arXiv:2307.13205*, 2023. 1
- [20] Mingcheng Li, Dingkang Yang, and Lihua Zhang. Towards robust multimodal sentiment analysis under uncertain signal missing. *IEEE Signal Processing Letters*, 2023. 2
- [21] Mingcheng Li, Dingkang 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)*, pages 10074–10082, 2024. 2
- [22] 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, 6, 7
- [23] 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, 6, 7, 8
- [24] 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
- [25] 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
- [26] 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 Con-*

- ference on Artificial Intelligence (AAAI), pages 2302–2310, 2021. 6, 7
- [27] 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)*, pages 5191–5198, 2020. 2
- [28] 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
- [29] 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. 2, 4
- [30] 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. 5
- [31] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. 5
- [32] 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)*, pages 6892–6899, 2019. 2, 6, 7, 8
- [33] 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. 2, 3, 4
- [34] Roece Shraga, Hagai 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
- [35] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 4
- [36] 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
- [37] 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, 6, 7
- [38] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019. 2, 4
- [39] 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. 3
- [40] 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. 2, 3
- [41] Lin Wang and Kuk-Jin Yoon. Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(6):3048–3068, 2021. 4
- [42] Shunli Wang, Dingkan Yang, Peng Zhai, Chixiao Chen, and Lihua Zhang. Tsa-net: Tube self-attention network for action quality assessment. In *Proceedings of the 29th ACM International Conference on Multimedia (ACM MM)*, pages 4902–4910, 2021. 1
- [43] Shunli Wang, Dingkan Yang, Peng Zhai, and Lihua Zhang. Cpr-clip: Multimodal pre-training for composite error recognition in cpr training. *IEEE Signal Processing Letters*, 2023. 1
- [44] 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)*, pages 7216–7223, 2019. 5
- [45] 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 (CVPR)*, pages 22025–22034, 2023. 2
- [46] Zilong Wang, Zhaohong Wan, and Xiaojun Wan. Transmodality: An end2end fusion method with transformer for multimodal sentiment analysis. In *Proceedings of The Web Conference*, pages 2514–2520, 2020. 2, 6, 7
- [47] 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. 2, 3
- [48] 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. 2
- [49] 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
- [50] 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. 2
- [51] 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.
- [52] 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.
- [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-founded emotion recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19005–19015, 2023. 1, 2
- [54] Dingkan Yang, Shuai Huang, Zhi Xu, Zhenpeng Li, Shunli Wang, Mingcheng Li, Yuzheng Wang, Yang Liu, Kun Yang, Zhaoyu Chen, Yan Wang, Jing Liu, Peixuan Zhang, Peng Zhai, and Lihua Zhang. Aide: A vision-driven multi-view, multi-modal, multi-tasking dataset for assistive driving perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20459–20470, 2023. 1
- [55] Dingkan 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, 2
- [56] 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. *arXiv preprint arXiv:2403.05023*, 2024.
- [57] Dingkan 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.
- [58] Dingkan Yang, Kun Yang, Mingcheng Li, Shunli Wang, Shuaibing Wang, and Lihua Zhang. Robust emotion recognition in context debiasing. In *CVPR*, 2024. 1, 2
- [59] Kun Yang, Dingkan Yang, Jingyu Zhang, Mingcheng Li, Yang Liu, Jing Liu, Hanqi Wang, Peng Sun, and Liang Song. Spatio-temporal domain awareness for multi-agent collaborative perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23383–23392, 2023. 1
- [60] Kun Yang, Dingkan Yang, Jingyu Zhang, Hanqi Wang, Peng Sun, and Liang Song. What2comm: Towards communication-efficient collaborative perception via feature decoupling. In *Proceedings of the 31th ACM International Conference on Multimedia (ACM MM)*, page 7686–7695, 2023. 1
- [61] 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. 2, 4
- [62] 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)*, pages 10790–10797, 2021. 1, 6, 7, 8
- [63] 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
- [64] 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. 5
- [65] 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. 5
- [66] 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
- [67] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11953–11962, 2022. 4
- [68] 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. 2
- [69] 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