Label-Guided Enhancement: A Distillation Framework for Uncertainty-Aware **Multimodal Emotion Recognition**

Abstract

1

Multimodal Emotion Recognition (MER) 2 is a vital technology for capturing nuanced 3 emotions human by integrating 4 complementary textual, acoustic, and 5 visual cues. However, real-world MER 6 systems frequently encounter issues with 7 missing or conflicting modalities, arising 8 from sensor failures, privacy constraints, or 9 contradictory emotional signals. These 10 issues compromise the efficacy of existing 11 attention-based fusion models. In this 12 propose CUMDF. paper, we а 13 Counterfactual-based Uncertain Missing 14 Modality **D**istillation **F**ramework that 15 addresses these challenges through three 16 core innovations. Specially, we introduce a 17 Label-Guided Multimodal Masked 18 Transformer (LG-MMT) to align features 19 with target sentiment semantics and 20 improving robustness under incomplete or 21 conflicting data. Furthermore, we design 22 the Adaptive and Generalized Knowledge 23 Extractors to disentangle modality-specific 24 information from shared cross-modal 25 patterns, enhancing representational 26 diversity and coherence. Finally, we design 27 Modality Attribution-based а 28 Counterfactual Inference (MACI) 29 mechanism that quantifies each modality's 30 causal contribution via counterfactual 31 predictions and dynamically adjusts 32 distillation weights to focus the student 33 model on under-optimized modalities. 34 Experimental results on three benchmark 35 demonstrate datasets that CUMDF 36 outperforms state-of-the-art approaches, 37 highlighting the importance of uncertainty 38 modeling in MER. 39

Introduction 40

43 signals (i.e. text, audio and visual) have become 44 indispensable for the accurate inference of 45 sentiment and emotion. Multimodal Emotion 46 Recognition (MER) employs a combination of 47 cues from these streams to enhance the accuracy of ⁴⁸ inference in comparison to unimodal systems, with applications in human-computer 49 extensive ⁵⁰ interaction. virtual assistants and affective ⁵¹ computing in the fields of healthcare and education. Despite the noteworthy advancements, two 52 ⁵³ under-explored challenges impede the robustness 54 and generalizability of existing MER systems. 55 Firstly, the absence of modalities is a common 56 occurrence in real-world implementations due to factors, including sensor 57 various failures, limitations, privacy 58 bandwidth filters. and 59 occlusion. This results in irregular availability of 60 modalities. Models that have been trained on fully 61 observed data encounter difficulties when 62 attempting to generalize under these incomplete 63 conditions. Second, modality conflicts—scenarios 64 where different modalities convey contradictory 65 sentiments (e.g., cheerful speech paired with a 66 frustrated facial expression)—are prevalent in 67 spontaneous human communication. It is evident 68 that standard fusion strategies are often insensitive 69 to such conflicts, resulting in semantically 70 incoherent feature representations and 71 compromised sentiment prediction.

To address these challenges, researchers have 72 73 explored various approaches to enhance MER 74 robustness. Early efforts focused on simple fusion 75 techniques using handcrafted features with 76 classical classifiers like SVMs (Rozgic et.al, 2012; 77 Wang et.al, 2014) and ensemble trees (Cummins 78 et.al, 2018). However, these methods proved ⁷⁹ incapable of capturing cross-modal dependencies, ⁸⁰ which are critical for emotion recognition (Smith et ⁸¹ al., 2017). The advent of deep learning gave rise to 41 With the explosive growth of user-generated 82 a more sophisticated array of fusion paradigms. ⁴² content on social media platforms, multimodal ⁸³ These include early fusion (Morency et al., 2011;

⁸⁴ Perez et al., 2013; Yu et al., 2021) which aggregates 123 refinements have been proposed, including 85 raw embeddings, and late fusion (Shutova et al., 124 margin-aware distillation (Wei et al., 2023) to ⁸⁶ 2016; Morvant et al., 2014; Evangelopoulos et al., 125 regularize modality importance, graph-based KD which combines 87 2013) 88 predictions. Nevertheless, these methodologies 127 information, 89 continue to encounter challenges in effectively 128 distillation (Wei et al., 2024) to improve teacher-⁹⁰ handling complex interactions between modalities. ¹²⁹ student alignment. Zhang et al. (2024) proposed a 91 ⁹² have become increasingly prominent. These ¹³¹ absent modality and optimization imbalance, ⁹³ include global attention modules (Zhang et al., ¹³² namely sample-weighted distillation and prototype 94 2020) and bi-GRU multi-attention (Lian et al., 133 regularization network. While these approaches 95 2021), which attempt to focus on salient 134 have advanced the field, they still suffer from three 96 multimodal features. Similarly, modality-weighted 135 fundamental limitations when dealing with real-97 fusion techniques, including gated networks 136 world MER scenarios: 98 (Arevalo et al., 2017) and self-adaptive path 137 ⁹⁹ selection (Yang et al., 2023), have been proposed ¹³⁸ generally perform cross-modal alignment without 100 to dynamically adjust each modality's contribution. 139 fine-grained label-aware supervision, rendering 101 Mai et al. (Mai et.al, 2023) introduced intra-modal 140 them susceptible to spurious correlations in noisy 102 and inter-modal comparative and 103 comparative learning. This instructional strategy is 142 anchoring to emotional content, models are unable 104 employed to facilitate a comprehensive exploration 143 to distinguish relevant features from noise, 105 of cross-modal interactions, to ensure the 144 particularly when modalities provide conflicting 106 maintenance of inter-class relationships, and to 145 signals. 107 address any existing modal gaps. Despite their 146 108 sophistication, these approaches typically assume 147 Current approaches fail to effectively separate 109 complete modality availability and ¹¹⁰ mechanisms to handle missing or conflicting ¹⁴⁹ information. This limits their ability to preserve 111 signals.

112 ¹¹³ missingness, knowledge distillation (KD) has been ¹⁵² robust performance when certain modalities are 114 identified as a promising approach. MCTN (Pham 153 missing or unreliable. et al., 2019) employs cycle-consistent translation 154 116 between modalities with a view to enhancing 155 methods apply identical distillation weights across ¹¹⁷ representation consistency across different input ¹⁵⁶ all modalities and samples, ignoring the varying 118 combinations. In alternative approaches, the 157 power and reliability of each modality. This ¹¹⁹ reconstruction of absent modalities is undertaken at ¹⁵⁸ becomes problematic 120 the feature level (Sun et al., 2023) or through the 159 conflicting modalities or different missing patterns 121 estimation of these modalities by means of 160 across samples. 122 completion-based KD (Sun et al., 2024). Further

modality-specific 126 (Deng et al., 2025) to aggregate multi-source and hierarchical cross-modal More recently, attention-based mechanisms 130 novel approach to address the challenges posed by

> Lack of semantic guidance: Existing methods semi- 141 modalities. In the absence of explicit semantic

Inadequate knowledge disentanglement: lack 148 modality-specific knowledge shared from 150 unique modality features while exploiting cross-With regard to the specific issue of modality 151 modality commonalities, which is critical for

> Consistent distillation strategies: Most when dealing with



Figure 1: The proposed CUMDF consists of three core components: (1) the Modality Domain 162

Knowledge Extractor (2) the Label-Guided Multimodal Masked Transformer, and (3) the Modality 163

164

Contribution-based Counterfactual Inference.

165 166 167 168 169 modalities. As demonstrated in Figure 1, CUMDF 215 modalities. introduces three key innovations: 170

17 Masked Transformer (LG-MMT), which employs 218 (MRMS), which generates datasets by dropping 172 173 learnable label embeddings as query inputs within 219 frame-level features with a drop ratio p (0 to 0.7, 174 the attention mechanism. In contrast to previous 220 0.1 increments). Unlike prior work (Yuan et.al, 175 fusion methodologies, the LG-MMT model 221 2021), this ensures at least one modality remains employs audio-text-label interactions with label 222 per sample for balanced evaluation, replacing embeddings as anchors to guide emotion feature 223 missing segments with zero vectors. Incomplete fusion. This approach focuses attention on 224 features X_m^s (for $m \in \{L, A, V\}$) are then processed 179 emotion-relevant features to stabilize cross-modal 225 by the student model similarly to the teacher model fusion and address issues related to missing or 226 (as shown in Eq. 1): 180 inconsistent signals. 181

Secondly, we propose the Adaptive Modality- 228 182 183 Specific (AMSKE) and the Generalized Modality- $_{229}$ sequence and d represents the feature dimension. 184 Common (GMCKE) Knowledge Extractors, which 185 decompose modality features into discriminative, 230 2.2 186 modality-specific traits and 187 and shared knowledge has been shown to enhance fusion, particularly in the presence of noisy or 190 absent streams. 191

Finally, the Modality Attribution-based 192 193 Counterfactual Inference (MACI) module adapts 194 distillation by estimating each modality's causal 239 195 impact via full-model and counterfactual 240 ¹⁹⁶ prediction comparisons, dynamically adjusting 197 distillation weights to promote fairness and 198 enhance under-optimized modal streams.

Methods 199 2

Preliminary 200 2.1

201 Multimodal Emotion Recognition (MER) is 202 commonly framed as a regression task. Given a ²⁰³ full-modality video sample set $S = [X_L, X_A, X_V]$, ²⁴⁹ where $X_L^t \in \mathbb{R}^{T_L * d_L}$, $X_V^t \in \mathbb{R}^{T_V * d_V}$, and $X_A^t \in \mathbb{P}^{250}$ ²⁰⁵ $\mathbb{R}^{T_A * d_A}$, respectively. And *t* represents the teacher ²⁵¹ 206 model, T_m is the sequence length (temporal ²⁵² 207 dimension) and d_m is the embedding dimension of ²⁵³ 208 modality $m \in \{L, A, V\}$. Meanwhile, we further ²⁵⁴ 209 extend the traditional MER task to encompass 255 \oplus denotes feature concatenation. The final output 210 modality-missing scenario, the modality missing

To address these limitations, the Counterfactual- 211 version features are denoted as X_m^s , where s based Uncertain Missing Modality Distillation 212 represents the student model, $m \in \{L, A, V\}$. Our Framework (CUMDF) is proposed, enhancing 213 goal is to find the utterance-level sentiments by MER performance under incomplete or conflicting 214 employing the multimodal data with missing

To handle intra-modality missingness, we 216 First, we propose a Label-Guided Multimodal 217 develop the Modality Random Missing Strategy

> $X'_{m}^{t} = Conv1D(X_{m}^{t}, k_{m}^{t}) + PE(T, d) \in \mathbb{R}^{T*d}$ (1) 227 where PE denotes to the position in the

Modality Knowledge Extractor (MKE)

cross-modal 231 To comprehensively represent multimodal inputs, commonalities. The explicit modelling of private 232 the framework uses Adaptive Modality-Specific 233 (AMSKE) and Generalized Modality-Common feature expressiveness and compatibility during 234 (GMCKE) Knowledge Extractors to decompose 235 unimodal features into modality-specific traits ²³⁶ $(X_m^t spe)$ and transferable cross-modal knowledge $_{237}$ ($X_{m \ com}^{t}$). Formally, the teacher model's extractors 238 are defined as:

3

$$X_{m_spe}^{t} = encoder_{m_spe} \left(X'_{m}^{t}; \theta_{m_spe} \right)$$
(2)

$$X_{m \ com}^{t} = encoder_{m \ com} \left(X_{m}^{\prime t}; \theta_{m \ com} \right)$$
(3)

where $encoder_{m_spe}$ and $encoder_{m_com}$ 241 ²⁴² denote the AMSKE and GMCKE, respectively. In 243 our framework, AMSKE employs a 1-layer 244 Transformer to learn sequential information of 245 modalities, and GMCKE is a fine-grained mamba ²⁴⁶ block. We use the language modality as an example 247 to present the structure of GMCKE, which is 248 calculated as follows:

$$M_{LA}^{t} = X_{L}^{\prime t} X_{A}^{\prime t^{T}}$$

$$S_{LA}^{t} = Softmax(tan(M_{LA}^{t}))$$

$$R_{LA}^{t} = S_{LA}^{t} X_{A}^{\prime t}$$

$$\tilde{X}_{LA}^{t} = R_{LA}^{t} \oplus (R_{LA}^{t} \odot X_{L}^{\prime t})$$

$$X_{L}^{t} com = \tilde{X}_{LA}^{t} \oplus \tilde{X}_{LV}^{t}$$

$$(4)$$

where \odot represents element-wise multiplication,

256 of the modality domain knowledge extractor $ilde{X}_m^t$ is 296 broadcast across the entire sequence, enabling the ²⁵⁷ a combination of $X_{m,spe}^{t}$ and $X_{m,com}^{t}$.

Label-guided Multimodal 258 2.3 Masked **Transformer (LG-MMT)** 259

261 often used for cross-modal relationship modeling. 301 into the Label-guided attention, where the label However, when modalities are absent or in 302 embedding e_i acts as query Q, \tilde{Z}_{m1}^i and \tilde{Z}_{m2}^i are 263 semantic conflict, Transformers may encounter 303 served as key K and value V. For instance, label-264 challenges due to an absence of global semantic 304 guided cross-attention between language and 265 guidance. To solve this, we propose the Label- 305 acoustic modalities focuses on features aligning 266 guided Multimodal Masked Transformer (LG- 306 with label sentiment, enhancing multimodal fusion: 267 MMT), as shown in Figure 2, which uses label 307 268 embeddings for each sample as semantic anchors 308 269 to guide attention during fusion.



271 Figure 2: Structure of the proposed LG-MMT

272 2.3.1 Label Embedding Generation

²⁷³ For each sample *i* with a label $y_i \in \{1, ..., C\}$, we ₂₇₄ map it into a dense vector $e_i \in \mathbb{R}^d$ using a ²⁷⁵ learnable label embedding matrix $E_{\nu} \in \mathbb{R}^{C*d}$:

 $e_i = E_v[y_i]$ 276

here, y_i denotes the sample label (discrete class ₃₂₃ 277 278 like "happy" or "sad"), with label embedding $e_i \in$ 279 \mathbb{R}^d encoding its semantic information. This 325 each attention head and passing them through a 280 embedding is crucial for guiding the attention 326 feed-forward neural network (FFN) for further 281 mechanism in subsequent layers.

282 2.3.2 Fusion of Label Embedding with Modality 283 Features

²⁸⁴ Let $\tilde{Z}_{m1}^i \in \mathbb{R}^d$ and $\tilde{Z}_{m2}^i \in \mathbb{R}^d$ be the feature 285 sequences of modalities m1and 286 m2 with m1, m2 $\in \{L, A, V\}$, respectively. The 287 label embedding e_i is fused with the modality 288 features to obtain the label-enhanced features. For 289 example, to obtain the label-enhanced language ²⁹⁰ feature \tilde{Z}_L^i , e_i is repeated along the time dimension 291 T_L to align with the temporal dimension of the ²⁹² language feature:

 $\tilde{Z}_{L}^{i} = \tilde{X}_{L}^{t} \odot (W_{e} * e_{i} + b_{e})$ (5)293 where W_e and b_e are learnable matrix and bias 294 295 vector. In this step, the label embedding is

297 model to use the label's semantic information ²⁹⁸ throughout the sequence.

299 2.3.3 Label-guided Attention Computation

260 In traditional MER methods, Transformers are 300 Label-enhanced features \tilde{Z}_{m1}^i and \tilde{Z}_{m2}^i are fed

$$Q = W_Q^{LE} * e_i$$

$$K_L = W_K^L * \widetilde{Z}_L^i$$

$$V_A = W_V^A * \widetilde{Z}_A^i$$
(6)

where W_0^{LE} , W_K^L and W_V^A are learnable weight 310 311 matrices for the query, key, and value 312 transformations, respectively. Next, we compute 313 the attention score and apply it to the value vectors ³¹⁴ to obtain the weighted features:

$$head_{LA} = Softmax \left(\frac{QK_L^{L}}{\sqrt{d_k}} + Mask\right) V_A \qquad (7)$$

$$Mask_m = \begin{cases} 0, & if token on the \\ padding position \\ 1, & otherwise \end{cases}$$
(8)

where Mask represents the mask matrix to 317 ³¹⁸ mask the attention tokens. By improving semantic 319 relevance between modalities and strengthening 320 intra-modality cohesion, the attention matrix's 321 weights more accurately reflect cross-modal (4) ₃₂₂ feature correlations.

After computing the attention heads, we 324 combine them by concatenating the outputs from 327 processing. This results in the final feature 328 representation for each modality:

$$H_{LA} = LayerNorm$$

$$\begin{pmatrix} \tilde{Z}_{L} + FFN(Concat(head_{LA}^{1}, \dots, head_{LA}^{h})W_{0}) \end{pmatrix} H^{s} = GAP(Concat(H_{LA}^{s}, H_{LV}^{s}, H_{AV}^{s}))$$
(9)
 $\hat{y}^{s} = Softmax(W_{s,pre}H^{s} + b_{s})$

Here, H_{LA}^{s} is the output of language and acoustic 332 333 features after attention and FFN processing; similarly, LG-MMT captures H_{LV}^s and H_{AV}^s . Global 335 Average Pooling (GAP) reduces fused feature $_{336}$ dimensionality to obtain H^s , which averages $_{337}$ features to retain critical information. Finally, H^s 338 feeds into a Softmax-connected fully connected ³³⁹ layer to compute prediction scores \hat{y}^s .

329

330

315

Similarity-based 340 2.4 Distillation 341

342 We introduce the Similarity-based Representation 389 ³⁴³ Distillation (SRD) to align the student and teacher 344 model representations, especially under missing 390 $_{345}$ modalities, via cosine similarity loss on fused $_{391}$ of modality m, these insights inform modality re-346 features. The similarity loss is defined as follows: 392 weighting,

$$L_{SRD} = \alpha_{st} \left(1 - \frac{W_s H^* \odot W_t H}{\parallel W_s H^t \parallel \parallel W_t H^s \parallel} \right)$$
(10)

348 349 hierarchical consistency and final representation 396 define the counterfactual loss as: $_{350}$ similarity, H^t and H^s represent the joint fused 351 representations from the teacher and student $_{352}$ models, respectively. The L_{SRD} penalizes any 353 divergence between H^t and H^s , encouraging the 354 student model to approximate the teacher's ³⁵⁵ multimodal representation as closely as possible.

356 2.5 Modality **Counterfactual Inference** 357

358 To quantify individual modality impacts on 404 matrices 359 predictions, we introduce the modality attribution- 405 Mathematically, it's defined as: based counterfactual inference (MACI) module. 406 360 By excluding specific modalities in counterfactual 407 407 $L_{reg} = \lambda_Q^{le} ||W_Q^{le}||^2 + \lambda_K^m ||W_K^m||^2 + \lambda_V^m ||W_V^m||^2$ (17) where W_Q^{le} , W_K^m , and W_V^m represent the weight 362 predictions, we assess their necessity for decision- $_{363}$ making. Formally, the removal of modality *m* from 364 the student model was shown to result in the 365 recompilation of output distributions. This process 366 defined uni-modal and counterfactual predictions 367 as follows:

 $\hat{y}_m^s = Softmax \left(W_s \tilde{X}_m^s + b_s \right)$ (11) 368 $\hat{y}_{con_m}^s = Softmax (W_{s_{con}} \tilde{X}_{w/o_m}^s + b_{s_{con}})$ (12) here, $\hat{y}_{con_m}^s$ is the counterfactual prediction 369 370 ³⁷¹ excluding modality m, with $\tilde{X}^{S}_{w/o_m} = \frac{1}{N} \sum_{i}^{N} \tilde{X}^{S}_{m}$ ⁴¹⁶ 372 as the student model's hidden representation $_{\rm 373}$ replacing m with its training-set average. The ³⁷³ replacing *m* with its training-set average. The ³⁷⁴ comparison of $\hat{y}_{con_m}^s$ with the full-modality ⁴¹⁸ sample, while \hat{y}_n represents the student model's ⁴¹⁹ predicted probability of the correct class. Finally, \hat{y}_m^s , enables modality *m*'s impact. $\frac{1}{420}$ we define the final loss function of CUMDF is ³⁷⁶ Finally, the attribution score, designated as δ_i^m , is $_{377}$ thus defined as the absolute difference between $_{422}$ 378 predictions $\hat{y}_{m(i)}^{s}$ and $\hat{y}_{con_{m}(i)}^{s}$: $\delta_i^m = \hat{y}_{m(i)}^s - \hat{y}_{con_m(i)}^s$ 379

where δ_i^m captures the predictive difference for 380 ³⁸¹ the *i-th* sample when modality m is omitted. Here, ⁴²⁵ **3 Experiments** $\hat{y}_{con_m(i)}^s$ represents the counterfactual prediction ⁴²⁶ 3.1 ³⁸³ of the student model, while $\hat{y}_{m(i)}^{s}$ corresponds to ³⁸⁴ the prediction of the student model under the same ⁴²⁷ We conduct extensive experiments on three MER 385 condition. To further quantify the relative 428 datasets with word-aligned data, including MOSI ³⁸⁶ importance of each modality across the dataset, we ⁴²⁹ (Zadeh et.al,2016), MOSEI (Zadeh et.al,2018), and

Representation 387 define the normalized modality contribution as 388 follows:

$$X_{con_m} = \frac{\sum_{i=1}^{N} |\delta_i^m|}{\sum_{m' \in \mathcal{M}} \sum_{i=1}^{N} |\delta_i^{m'}|}$$
(14)

where X_{con_m} denotes the attribution proportion feature selection. and model ³⁹³ compression in resource-constrained settings. To ³⁹⁴ ensure the model learns modal relative importance where α_{st} is the weighting vector to balance 395 while maintaining prediction consistency, we

$$\lambda_{con_m} = \frac{X_{con_m}}{\sum_{m' \in M} X_{con_{m'}}} \tag{15}$$

$$L_{cou} = \lambda_{con_m} \sum_{m \in M} X_{con_m}$$
(16)

where λ_{con_m} is a normalization coefficient that 399 400 ensures the modality attributions sum to 1.

To avoid overfitting and promote robust learning, 401 Attribution-based $_{402}$ we use a regularization loss L_{reg} that penalizes the 403 magnitudes of the attention mechanism's weight (query, key, value projection).

> 408 matrices for the query, key, and value 409 transformations in the attention layers of the model, 410 while λ_Q^{le} , λ_K^m , and λ_V^m denote their corresponding 411 regularization coefficients. Finally, we defined the 412 task-specific loss L_{task} to directly optimize the 413 model's performance on the sentiment prediction 414 task We utilize the conventional cross-entropy loss, 415 which is characterized as:

$$y_{task} = -\frac{1}{N} \sum_{n=1}^{N} y_n \log \hat{y}_n$$
 (18)

where y_n denotes the true label for the *n*-th 421 defined as Eq. (19):

 $L = \lambda_1 L_{reg} + \lambda_2 L_{cou} + \lambda_3 L_{SRD} + L_{task}$ (19)where λ_1 , λ_2 , and λ_3 are weights that control 423 (13) $_{424}$ the relative importance of each loss term.

Datasets and Evaluation Metrics

431 dataset are shown in Appendix A.1.

432 3.2 **Implementation Details**

434 framework, utilizing NVIDIA Tesla V100 GPUs 445 TransM (Wang et.al, 2020). To simulate missing 435 and torch version 1.8.2. The details of parameter 446 inter-modalities, we remove modalities from ⁴³⁶ implementations are listed in Appendix A.2.

Comparison with the state-of-the-art 437 3.3

438 We compare our CUMDF with eight exemplary

439 and replicable state-of-the-art (SOTA) approaches,

Table 1: Performance comparison results on MOSI and MOSEI

	Models	Testing Conditions													
Dataset		{l}		{a}		{v}		{ <i>l</i> , <i>a</i> }		{ <i>l</i> , <i>v</i> }		$\{a, v\}$		$\{l, a, v\}$	
		MAE	F1	MAE	F1	MAE	F1	MAE	F1	MAE	F1	MAE	F1	MAE	F1
MOSI	Self-MM	0.810	66.25	0.764	41.37	0.752	39.15	0.753	68.94	0.742	73.89	0.852	47.90	0.814	83.15
	CubeMLP	0.804	65.78	0.782	43.10	0.748	41.20	0.801	64.76	0.695	68.12	0.823	49.04	0.803	79.82
	DMD	0.795	67.45	0.801	43.65	0.761	42.90	0.698	69.42	0.735	67.93	0.831	51.25	0.892	82.14
	MCTN	0.821	75.56	0.715	58.96	0.797	58.12	0.710	76.43	0.703	73.62	0.864	62.47	0.802	83.75
	EMT	0.845	63.87	0.807	39.25	0.739	43.47	0.721	63.15	0.734	65.85	0.845	48.36	0.789	83.67
	IF-MMIN	0.830	56.14	0.841	46.89	0.815	45.13	0.745	61.27	0.785	64.94	0.802	66.92	0.762	81.24
	TransM	0.857	58.21	0.851	66.83	0.823	51.09	0.749	58.45	0.739	68.92	0.804	65.37	0.785	82.67
	CUMDF	0.766	81.93	0.745	62.29	0.726	60.48	0.689	76.38	0.675	80.41	0.731	75.29	0.718	84.37
MOSEI	Self-MM	0.763	67,85	0.731	43.83	0.725	41.74	0.775	69.43	0.688	72.52	0.803	48.95	0.752	82.74
	CubeMLP	0.772	71.65	0.749	43.52	0.736	37.67	0.722	75.94	0.725	74.67	0.824	49.58	0.739	83.07
	DMD	0.796	70.37	0.784	42.37	0.755	38.42	0.684	74.55	0.713	72.68	0.817	50.26	0.807	81.75
	MCTN	0.809	72.93	0.695	44.26	0.788	39.18	0.692	74.21	0.691	73.63	0.848	62.91	0.728	82.78
	EMT	0.791	67.46	0.716	39.55	0.719	32.57	0.700	71.63	0.708	70.05	0.822	48.52	0.780	83.13
	IF-MMIN	0.823	68.59	0.758	41.47	0.781	33.96	0.719	71.26	0.752	70.73	0.796	48.95	0.743	82.76
	TransM	0.801	69.75	0.782	42.56	0.803	34.98	0.727	73.28	0.729	71.19	0.773	50.34	0.775	82.91
	CUMDF	0.742	72.48	0.735	57.53	0.717	58.45	0.678	79.72	0.676	81.03	0.735	71.57	0.704	83.25

450 451

449

As illustrated in Table 1, CUMDF demonstrates 466 a feature that is particularly advantageous when 452 efficacy in modality missing scenarios in 467 language is missing.

453 comparison to baselines for both MOSI and 468 454 MOSEI. In the context of MOSI in the language 469 performance of the CUMDF and baseline models 455 scenario, CUMDF exhibits the lowest MAE of 470 is shown for various modality circumstances on 456 0.766 and the highest F1 of 81.93%. This 471 IEMOCAP. In the "Happy" emotion category, the 457 performance surpasses the second-best by 6.37% in 472 CUMDF achieves the highest F1 scores in three out 458 terms of F1. This phenomenon can be attributed to 473 of seven testing conditions, including single-459 the alignment of features with label sentiment 474 modality scenarios {1} (82.4%), {a} (69.3%), and 460 priors by LG-MMT, a process that ensures the 475 {V} (68.1%). Similarly, the "Sad" category extraction of robust information. 461

462 463 75.29%, which is 9.92% higher than the next best. 478 The adaptability of the model's knowledge 464 This demonstrates the efficacy of MACI, which 479 extraction and fusion mechanisms is emphasized 465 adjusts modal contributions based on causal impact,

As demonstrated in Table 2, the F1 score 476 exhibited consistent superiority of our method over In the {a,v} scenario on MOSI, CUMDF's F1 is 477 baselines across five distinct testing conditions.

430 IEMOCAP (Busso et.al, 2008). The statistics of the 440 including complete-modality methods: Self-MM 441 (Yu et.al, 2021), CubeMLP (Sun et.al, 2022), and 442 DMD (Li et.al, 2023) and missing-modality 443 methods: MCTN (Pham et.al, 2019), EMT (Sun 433 All models are constructed using the Pytorch 444 et.al, 2022) IF-MMIN (Zuo et.al, 2023) and 447 samples in Table 1 and Table 2. "{1}" means only 448 language is present.

480 by its consistent superiority across diverse emotion

481 categories.

	Table 2	: FI score	performan	ice compa	rison resu	lts on IEMC	DCAP	
Models	Matria	Testing Conditions						
widueis	wienies	$\{l\}$	{ <i>a</i> }	{ <i>v</i> }	{ <i>l</i> , <i>a</i> }	$\{l, v\}$	{a, v}	$\{l, a, v\}$
	Нарру	68.1	53.4	51.3	72.5	69.7	62.0	88.6
Salf MM	Sad	69.2	53.0	53.8	68.9	68.1	61.8	87.9
Sen-Iviivi	Angry	67.5	54.2	53.3	67.1	68.9	57.2	85.7
	Neutral	56.2	48.5	50.7	58.3	56.2	52.6	69.8
	Нарру	66.2	52.3	50.1	69.4	68.5	56.2	88.9
CuboMI D	Sad	68.5	51.8	54.7	71.8	69.5	57.5	86.7
Cubewill	Angry	65.2	53.7	51.6	69.5	67.8	56.6	85.4
	Neutral	55.7	48.6	50.3	58.5	56.9	52.8	70.7
	Нарру	69.8	55.7	52.2	79.3	77.5	65.8	83.1
DMD	Sad	65.5	55.7	53.6	78.3	73.4	68.7	82.8
DND	Angry	65.3	54.1	51.4	81.7	80.4	59.5	84.6
	Neutral	54.6	51.3	49.5	65.3	62.7	54.9	67.1
	Нарру	77.5	63.2	61.7	81.3	80.4	66.5	85.5
MTCN	Sad	76.2	64.3	60.5	82.9	81.5	64.3	84.0
MICN	Angry	77.1	61.5	58.4	83.7	81.5	68.4	85.1
	Neutral	60.5	51.2	50.6	65.8	62.7	56.5	67.1
	Нарру	68.3	54.2	51.5	72.1	69.4	60.3	89.0
EMT	Sad	65.7	54.6	53.4	70.5	68.73	58.1	88.5
	Angry	65.2	53.6	50.9	69.7	69.5	54.8	86.1
	Neutral	53.8	50.4	48.5	57.2	54.3	51.8	71.8
	Нарру	80.3	66.8	64.5	83.4	81.7	67.2	90.3
IF MMIN	Sad	79.2	65.9	62.3	82.6	79.5	70.4	85.2
	Angry	80.1	67.6	61.2	83.4	82.3	59.7	84.9
	Neutral	61.2	50.8	49.5	62.4	52.9	55.7	67.2
	Нарру	82.3	67.7	66.9	83.5	82.6	69.8	87.3
TurneM	Sad	81.7	69.5	66.3	84.1	81.8	70.3	86.9
	Angry	81.5	67.6	65.5	82.5	81.6	68.1	85.2
	Neutral	61.2	52.3	43.1	64.9	62.7	57.2	71.5
	Нарру	82.4	69.3	68.1	84.2	82.1	70.1	87.6
CUMDE	Sad	82.5	71.7	67.5	83.3	82.4	72.4	88.7
CUMDF	Angry	82.7	67.4	66.3	83.6	82.9	67.3	86.2
	Neutral	63.3	54.1	52.4	68.3	64.5	57.8	71.2

482

. .

483

To illustrate performance variation at different 488 scores for the "happy" and "sad" categories. This is 484 485 drop rates, Figure 3 plots the performance curves 489 due to its AMSKE and GMCKE modules, which 486 of models across intra-modality drop ratios (0.1- 490 extract modality-specific and modality-common 487 0.7). CUMDF maintains significantly higher F1 491 knowledge to compensate for missing information.



Figure 3: F1-score results of different intra-modality drop ratios on IMEOCAP.

494 3.4 Ablation Studies of Modules

496 our framework, we conducted comprehensive 505 to measure modal causal impacts and adaptively 497 ablation studies on the MOSI dataset with a drop 506 weight them during fusion. Removing MKE leads ratio of 0.2, as shown in Table 3.

499 500 causes significant performance drops, confirming 509 is crucial for capturing unique traits and cross-501 the label-guided attention mechanism focuses on 510 modal commonalities.

502 emotion-relevant features and improves cross-503 modal alignment. Removing MACI also degrades 495 To validate the contribution of each component in 504 performance, validating counterfactual reasoning 507 to severe performance deterioration, indicating Ablation studies show removing LG-MMT 508 disentangling modality-specific/shared knowledge

511	Table 3: Module ablation performance results on CMU-MOSI with a drop ratio of $p = 0.2$.						
	MKE	LG-MMT	MACI	MAE	F1		
				0.820	82.36		
		\checkmark	\checkmark	0.782	84.15		
	√		\checkmark	0.798	83.82		
	√	\checkmark		0.764	84.09		
	√	√	√	0.758	84.57		

512

513 3.5 **Case Study**

Figure 4 shows CUMDF versus baseline IF- 519 prioritize reliable 514 515 MMIN on three MOSI test samples, with missing 520 predictions closer to ground truth labels. 516 modalities marked by grey rectangles. CUMDF

517 outperforms the baseline across missing modality 518 scenarios by using counterfactual reasoning to modal signals, vielding



522

Figure 4: The case of emotion recognition base on our proposed CUMDF.

Conclusion 523

524 525 distillation framework for multimodal sentiment 541 CUMDF consistently outperforms state-of-the-art 526 analysis that effectively addresses modality 542 approaches across various modality missing missingness and conflicts. We design a Label- 543 scenarios. 527 Guided Multimodal Masked Transformer (LG-MMT) that incorporates label embeddings as 544 529 530 semantic anchors to guide cross-modal feature ⁵³¹ fusion. The proposed Adaptive Modality-Specific and Generalized Modality-Common Knowledge Extractors disentangle unique modality traits from 534 shared patterns, enhancing representational 535 coherence. Additionally, our Modality Attribution-536 based Counterfactual Inference 537 mechanism dynamically adjusts

538 contributions based on their causal impact on 539 prediction. Extensive experiments on MOSI, In this paper, we introduce CUMDF, a novel 540 MOSEI, and IEMOCAP datasets demonstrate that

Limitation

Our CUMDF has the following limitations: (1) Multimodal Masked 546 The Label-Guided 547 Transformer and counterfactual inference 548 mechanisms increase computational complexity 549 compared to simpler fusion approaches. (2) The (MACI) 550 effectiveness of CUMDF depends on the quality of modality 551 initial modality representations and could benefit 552 from more advanced feature encoders, though this ⁵⁵³ is not the focus of our work. (3) While we evaluated ⁶⁰³ 554 our framework on standard EMR benchmarks, it 604 ⁵⁵⁵ lacks validation on more diverse scenarios such as 556 real-time streaming data with dynamic modality 606 Verónica Pérez-Rosas, Rada Mihalcea, and Louis-557 availability or datasets from specialized domains 607 558 like healthcare or education. Testing across 608 559 additional multimodal tasks beyond emotion 609 610 ⁵⁶⁰ recognition could further validate the framework's ⁵⁶¹ generalizability and effectiveness in broader ⁶¹¹ Wenmeng Yu, Hua Xu, Ziqi Yuan, and Jiele Wu. 2021. 612 562 contexts.

Acknowledgments 563

564 This work was supported by the National Social 616 565 Science Fund of China (22BTQ048), Scientific 617 Ekaterina Shutova, Douwe Kiela, and Jean Maillard. 566 Research Project of Education Department of Jilin 618 567 Province (JJKH20250758KJ), and the Projects of 619 568 Jilin University of Finance and Economics 620 621 569 (2023YB021 and 2024PY010). 622

570 References

571 Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-625 Philippe Morency, and Barnabás Póczos. 2019. 572 626 Found in translation: Learning robust joint 573 627 representations by cyclic translations between 574 628 modalities. In Proceedings of the AAAI Conference 575 on Artificial Intelligence, pages 6892-6899. 576

577 Licai Sun, Zheng Lian, Bin Liu, Jianhua Tao, 2023. Efficient multimodal transformer with dual-level 578 feature restoration for robust multimodal emotion 579 recognition. IEEE Transactions on Affective 580 Computing. 581

iktor Rozgić, Sankaranarayanan Ananthakrishnan, 582 Shirin Saleem, Rohit Kumar, Rohit Prasad. 2012. 583 Ensemble of SVM trees for multimodal emotion 584 recognition. In Proceedings of the 2012 Asia Pacific 585 Signal and Information Processing Association 586 Annual Summit and Conference, pages 1-4. 587 588 Nicholas Cummins, Shahin Amiriparian, Sandra Ottl,

Maurice Gerczuk, Maximilian Schmitt, Björn 589 Schuller. 2018. Multimodal bag-of-words for cross 590 domains emotion recognition. In Proceedings of the 591 2018 IEEE International Conference on Acoustics, 592 Speech and Signal Processing (ICASSP), pages 593 4954-4958. 594 595 Min Wang, Donglin Cao, Lingxiao Li, Shaozi Li, and

649 Rongrong Ji. 2014. Microblog Emotion recognition 596 650 Based on Cross-Media Bag-of-Words Model. In 597 Proceedings of International Conference on 598 Internet Multimedia Computing and Service 652 Zheng Lian, Bin Liu, and Jianhua Tao. 2021. CTNet: 599 (ICIMCS), pages 76-80. 600 653

601 Louis-Philippe Morency, Rada Mihalcea, and Paresh Doshi. 2011. Towards multimodal emotion 602

recognition: Harvesting opinions from the web. In Proceedings of the 13th International Conference on Multimodal Interfaces (ICMI), pages 169-176.

Philippe Morency. 2013. Utterance-level multimodal emotion recognition. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, pages 973–982.

Learning modality-specific representations with self-supervised multi-task learning for multimodal 613 emotion recognition. In Proceedings of the Thirty-614 Fifth AAAI Conference on Artificial Intelligence, 615 pages 10790-10797.

2016. Black holes and white rabbits: Metaphor identification with visual features. In, The 2016 Conference of the North Ameri can Chapter of the Association for Computational Linguistics (NAACL) HLT 2016: Human Language Technologies, pages 160-170.

- 624 Emilie Morvant, Amaury Habrard, and Stéphane Ayache. 2014. Majority vote of diverse classifiers for late fusion. Structural, Syntactic, and Statistical Pattern Recognition - In Joint IAPR International Workshop, S+SSPR 2014, pages 153-162. Springer.
- 629 Georgios Evangelopoulos, Athanasia Zlatintsi. Alexandros Potamianos, Konstantinos Rapantzikos, 630 and Georgios Skoumas. 2013. Multimodal saliency 631 and fusion for movie summarization based on aural, visual, and textual attention. IEEE Transactions on 633 Multimedia, 15(7): 1553-1568. 634
- 635 John Arevalo, Thamar Solorio, Manuel Montes-y-Gómez, and Fabio A. González. 2017. Gated 636 multimodal units for information fusion. In 638 Proceedings of the International Conference on Learning Representations, pages 1–17. 639
- 640 Haozhe Yang, Xianqiang Gao, Jianlong Wu, Tian Gan, Ning Ding, Feijun Jiang, and Liqiang Nie. 2023. 641 Self-adaptive Context and Modal-interaction 642 Modeling for Multimodal Emotion Recognition. In 643 Proceedings of the 61th annual meeting of the 644 Association for Computational Linguistics, pages 645 6267-6281. 646
- 647 Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, et al. 2020. Relation-aware global attention for person re-648 identification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2020), pages 3183-3192.
- Conversational transformer network for emotion recognition. IEEE/ACM Transactions on Audio 654 Speech and Language Processing, 29: 985–1000.

656 Sijie Mai, Ying Zeng, Shuangjia Zheng, and Haifeng 710 Hao Sun, Hongyi Wang, Jiaqing Liu, Yen-Wei Chen,

Hu. 2023. Hybrid contrastive learning of tri-modal 711 657

658 representation for multimodal emotion recognition. 712

IEEE Transactions on Affective Computing, 14(3): 713 659

714

719

730

745

2276-2289. 660

715 661 Yuhang Sun, Zhizhong Liu, Quan Z. Sheng, Dianhui Chu, Jian Yu, and Hongxiang Sun. 2024. Similar 716 Yong Li, Yuanzhi Wang, and Zhen Cui. 2023. 662 modality completion-based multimodal emotion 717 663 recognition under uncertain missing modalities. 718 664

Information Fusion, 110: 114. 665

720 666 Shicai Wei, Yang Luo, and Chunbo Luo. 2023. 667 aware regularization for incomplete multimodal 722 668 learning. In 2023 IEEE/CVF Conference on 723 669 Computer Vision and Pattern Recognition, pages 724 670 20039-20049. 671 725

672 Yuanyue Deng, Jintang Bian, Shisong Wu, Jianhuang 726 Hamlin Zuo, Rui Liu, Jinming Zhao, Guanglai Gao, Lai, and Xiaohua Xie. 2025. Multiplex graph 727 673

- aggregation and feature refinement for unsupervised 728 674
- incomplete multimodal emotion recognition. 729 675

Information Fusion, 114: 102711. 676

731 677 Puling Wei, Juan Yang, and Yali Xiao. 2024.

678

- Network Enhanced with Self-Distillation for 733 679 Emotion Recognition in Conversations. Electronics, 734 680
- 13:2645. 681

736 Yujuan Zhang, Fang'ai Liu, Xuqiang Zhuang, Ying 682 683 sample-weighted distillation unified framework 738 684 adapted to missing modality emotion recognition. 739 685 Neural Networks, 177: 106397. 686 740

741 687 Ziqi Yuan, Wei Li, Hua Xu, and Wenmeng Yu. 2021. 742 Transformer-based feature reconstruction network 688 for robust multimodal emotion recognition. In 743 W 689

Proceedings of the 2021 ACM Multimedia 744 690

- Conference (MM '21), pages 4400-4407. 691
- 746 692 Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-747 Philippe Morency. 2016. MOSI: Multimodal corpus 693 748 of sentiment intensity and subjectivity analysis in 694 online opinion videos. arXiv 695
- arXiv:1606.06259. 750 696

751 697 Amir Zadeh, Paul Pu Liang, Soujanya Poria, Erik 752 Cambria, and Louis-Philippe Morency. 2018. 698 Multimodal language analysis in the wild: CMU-699 700 graph. In Proceedings of the 56th Annual Meeting of 755 701

- the Association for Computational Linguistics, 756 702
- pages 2236-2246. 703

758 Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe 704 759 Kazemzadeh, Emily Mower, Samuel Kim, Jeannette 705 706 Narayanan. 2008. IEMOCAP: Interactive emotional 761 707 dyadic motion capture database. Language 762 708 Resources and Evaluation, 42: 335-359. 709

- and Lanfen Lin. 2022. CubeMLP: An MLP-based Model for Multimodal Emotion recognition and Depression Estimation. In Proceedings of the 30th ACM International Conference on Multimedia (MM '22), pages 3722-3729.
- Decoupled multimodal distilling for emotion recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6631-6640.
- MMANet: Margin-aware distillation and modality- 721 Sun Licai, Lian Zheng, Liu Bin, and Tao Jianhua. 2022. Efficient Multimodal Transformer with Dual-Level Feature Restoration for Robust Multimodal Emotion recognition. arXiv preprint arXiv:2208.07589.
 - and Haizhou Li. 2023. Exploiting modalityinvariant feature for robust multimodal emotion recognition with missing modalities. In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1-5.
- Hierarchical Cross-Modal Interaction and Fusion 732 Zilong Wang, Zhaohong Wan, and Xiaojun Wan. 2020. Transmodality: An end2end fusion method with transformer for multimodal emotion recognition. In Proceedings of The Web Conference, pages 2514-735 2520.
- Hou, and Yuling Zhang. 2024. Prototype-based 737 Yansen Wang, Ying Shen, Zhun Liu, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. 2019. Words can shift: Dynamically adjusting word representations using nonverbal behaviors. In Proceedings of the AAAI Conference on Artificial Intelligence, pages 7216–7223.
 - enbin Wang, Liang Ding, Li Shen, Yong Luo, and Zhe Li. 2023. Transformer-based multimodal feature fusion for emotion recognition with missing modalities. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8004-8012.
 - preprint 749 Jeffery Pennnington, Richard Socher and Christopher D. Manning. 2014. GloVe: Global vectors for word representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. Pages 1532–1543.
- MOSEI dataset and interpretable dynamic fusion 754 Gilles Degottex, John Kane, Thomas Drugman, Tuomo Raitio and Stefan Scherer. 2014. COVAREP: A collaborative voice analysis repository for speech technologies. In IEEE International Conference on 757 Acoustics, Speech and Signal Processing, pages 1532-1543.
- N. Chang, Sungbok Lee, and Shrikanth S. 760 Zengqun Zhao, Qingshan Liu and Shanmin Wang (2021) Learning deep global multi-scale and local attention features for facial expression recognition

in the wild. IEEE Trans Image Process 30:6544- 810 A.2 Implementation Details 763 6556. 764

Penn Phonetics Laboratory (2013) p2fa-vislab: A script 765 F

766

from: https://github.com/ucbvislab/p2favislab/ 767

769 770 771 772 pages 1545-1554. 773

Appendices 774 A

775 A.1 Datasets and Evaluation Metrics

We conduct extensive experiments on three 776 777 MER datasets with word-aligned data, including MOSI, MOSEI, and IEMOCAP. The MOSI dataset consists of 2,199 video clips containing opinion 779 comments from 89 independent speakers debating 780 YouTube movie reviews. The sample comprises 41 782 female and 48 male speakers. Every video clip is 783 labelled with sentiment strength, from -3 (very 784 negative) to +3 (extremely positive). The MOSEI 785 is a dataset consisting of 22,856 video clips, which 786 has 16,326, 1,871, and 4,659 samples in train, valid, 834 attributes 787 and test data. Each sample of MOSI and MOSEI is 788 labeled by human annotators with a sentiment 789 score of -3 (strongly negative) to +3 (strongly 790 positive).

On the MOSI and MOSEI datasets, we utilize 791 792 weighted F1 score computed for positive/negative classification results as evaluation metrics. The 794 IEMOCAP dataset is intended for multi-label 795 emotion recognition. It comprises 302 videos, ⁷⁹⁶ including 151 recorded chat videos. Every sentence ⁷⁹⁷ in the dialogue snippets is categorized by a specific 798 emotion: happiness, sadness, anger, surprise, fear, 799 and ten additional emotions. In our research, we follow the idea from (Wang et.al, 2019) and 800 concentrate on identifying four fundamental 801 emotions: happiness, sadness, anger, and neutrality. 802 According to a previous experimental study 846 All models are constructed using the Pytorch 803 804 (Wang et.al, 2023), our model's performance 847 framework, utilizing NVIDIA Tesla V100 GPUs 805 leverages two metrics, adopting a dual approach for 848 and torch version 1.8.2. To facilitate a fair 807 regression, we provide MAE as a measurement. 850 (SOTA) approaches utilizing publicly accessible 808 For the classification tasks, we provide accuracy 851 codebases 809 F1-score (F1) as measurement.

811 A.2.1 Feature Extraction

for audio/transcript alignment. GitHub. Available 812 For the language modality, we transform the video 813 transcripts into pre-trained GloVe (Pennnington 768 Jiandian Zeng, Tianyi Liu and Jiantao Zhou. 2022. Tag- 814 et.al, 2014) word embeddings to acquire a 300assisted multimodal emotion recognition under 815 dimensional vector. In the audio modality, we uncertain missing modalities. In Proceedings of the 816 employ the COVAREP (Degottex et.al, 2014) to 45th International ACM SIGIR Conference on 817 extract 74-dimensional low-level audio features, Research and Development in Information Retrieval, 818 which can process Mel-cepstral coefficients, 819 fundamental frequency, voiced/unvoiced segments, 820 normalized amplitude quotient, quasi-open 821 quotient, glottal source parameters, harmonic 822 model, phase distortions, and formants. MA-Net 823 (Zhang et.al, 2021) is utilized as the video feature 824 extractor for the video modality. Renowned for its 825 considerable facial success in expression 826 recognition, we initially employed the MTCNN 827 face detection algorithm to identify faces. 828 Subsequently, we use the pre-trained MA-Net ⁸²⁹ model to extract 1024-dimensional video features.

> To attain word-level alignment among the three 830 ⁸³¹ modalities, we first process the video and audio 832 streams via P2FA (Penn, 2013) to provide aligned 833 timestamps. Thereafter, the video and audio are averaged throughout these 835 synchronized intervals. In the CMU-MOSI dataset, 836 the sequence lengths for all three modalities are 837 established at 50. Conversely, the other two dataset 838 maintain a sequence length of 20 throughout all 839 three modalities.

> 840 Additionally, we partition these two datasets into ⁸⁴¹ training sets, validation sets, and test sets according 842 to established ratios, with the dataset sizes 843 specified in Table 1.

Table 1: Statistics of Datasets

Dataset	CMU-MOSI	CMU-MOSEI	IEMOCAP
Training Set	1284	16326	2717
Validation Set	229	1871	798
Test set	686	4659	938

845 A.2.2 Experimental Setup

classification and regression predictions. For 849 comparison, we re-implement the state-of-the-art and integrate them with our frameworks. All 852 experimental experimental 853 findings are averaged across numerous trials with 854 five distinct random seeds.

The hyperparameters of our proposed model utilize the configurations outlined in (Zeng et.al,2022) as detailed in Table 2. The learning rate ses is established as 0.001, and the concealed size is designated as 300. We employed the Adam optimizer to reduce the overall loss function and enhance the proposed AUMDF. The duration is established at 20, and the loss weight is designated at 0.1. The batch size for the CMU-MOSI dataset is set at 35, while the other two are designated as 32.

866 Table 2: Hyperparameter settings of CUMDF

rable 2. Hyperparameter settings of COMDI						
Description	Symbol	Value				
Epoch	b	20				
Dropout rate	d	0.8				
Hidden size	h	300				
Learning rate	lr	0.001				
Drop ratio	p	[0.1-0.7]				
Maximum language length	m_l	25				
Maximum audio length	m_a	150				
Maximum video length	m_v	100				
Loss weights	$\lambda_1, \lambda_2, \lambda_3$	0.1				
Early stop	es	20				