

# Incongruity-Aware Multimodal Sentiment and Emotion Recognition

Anonymous EMNLP submission

## Abstract

Fusing multiple modalities for affective computing tasks has proven effective for performance improvement. However, how multimodal fusion works and how to use it are not well understood. Also, multimodal fusion is difficult to use in real-world scenarios due to large model sizes. In this work, on sentiment and emotion recognition, we first analyze how the salient affective information in one modality can be affected by the other in crossmodal attention. We find that inter-modal incongruity exists at the latent level due to crossmodal attention. Based on this finding, we propose a lightweight model via Hierarchical Crossmodal Transformer with Modality Gating (HCT-MG), which determines a primary modality according to its contribution to the target task and then hierarchically incorporates auxiliary modalities to alleviate inter-modal incongruity and reduce information redundancy. The experimental evaluation on three benchmark datasets: CMU-MOSI, CMU-MOSEI, and IEMOCAP verifies the efficacy of our approach, showing that it: 1) outperforms major prior work by achieving competitive results and recognizing hard samples; 2) mitigates the inter-modal incongruity at the latent level when modalities have mismatched affective tendencies; 3) reduces model size to less than 1M parameters while outperforming existing models of similar sizes.

## 1 Introduction

As emotions are expressed in complex ways (e.g., face, voice, and language) in human communication, multimodal fusion has become a hot topic in the past decade. Previous studies have shown that by taking advantage of complementary information from multiple modalities, emotion recognition can be more robust and accurate (Xu et al., 2018; Li et al., 2022a). However, several major issues remain unsolved, impeding the progress of Multimodal Sentiment and Emotion Recognition

(MSER). First, multimodal signals often show an unaligned nature, bringing about the asynchrony problem (Tsai et al., 2019). For example, the visual signal usually precedes the audio by around 120ms when people express emotion (Grant and Greenberg, 2001). Second, different modalities may have different or even opposite affective tendencies, which makes emotions difficult to recognize. For example, people sometimes say positive content with a negative voice (e.g., sarcasm) or negative content with a smile (e.g., to be polite).

Prior work has proposed many approaches to tackle these issues. For example, Tsai et al. (2019) introduced the Multimodal Transformer (MulT) model to learn a pair-wise latent alignment with the Transformer structure, which directly attends to low-level features in multiple modalities to solve the asynchrony problem. Wu et al. (2021) proposed an incongruity-aware attention network that focuses on the word-level incongruity between modalities by assigning larger weights to words with incongruent modalities. Nevertheless, to capture as much information as possible for better performance, recent models usually repeatedly fuse specific or all modalities (Liang et al., 2018), resulting in not only redundant information but also large model sizes that hinder their real-world use.

To address this problem, in this paper we propose the Hierarchical Crossmodal Transformer with Modality Gating (HCT-MG), a lightweight multimodal fusion model that can alleviate inter-modal incongruity, reduce information redundancy, and learn representations from unaligned modalities at the same time. Specifically, HCT-MG dynamically determines the primary modality based on its contribution to the target task and then hierarchically fuses auxiliary modalities via crossmodal Transformers to efficiently obtain the most useful information without modality alignment. The model is motivated by a feasibility analysis of the crossmodal Transformer (specifically, its attention

mechanism) in multimodal fusion (i.e., how the salient affective information in one modality is affected by the other at the latent level).

The feasibility analysis demonstrates that cross-modal attention functions by highlighting the salient affective information in one modality with the help of the other one. However, when modalities have mismatched affective tendencies, cross-modal attention may malfunction by leaving inter-modal incongruity at the latent level. The experimental evaluations on CMU-MOSI (Zadeh et al., 2016), CMU-MOSEI (Zadeh et al., 2018b), and IEMOCAP (Busso et al., 2008) show that our approach achieves competitive results and alleviates the inter-modal incongruity with a small model size. (Code available<sup>1</sup>)

## 2 Related Work

Among previous approaches, early fusion and late fusion are the most widely used for MSER. However, due to the strict constraint on time synchrony, early fusion does not work well if the input features of multiple modalities differ in their temporal characteristics (Li et al., 2020). On the other hand, since different modalities have been confirmed to be complementary to each other (Chuang and Wu, 2004), the relatedness among them is ignored by late fusion. To this end, tensor fusion, which is performed at the latent level, has become mainstream. For example, Zadeh et al. (2017) introduced a Tensor Fusion Network, that learns both intra- and inter-modality dynamics end-to-end.

Furthermore, with the success of the cross-attention mechanism (Lu et al., 2019), which exchanges key-value pairs in self-attention, a major trend using cross-attention for multimodal fusion has emerged and is usually referred to as *crossmodal attention*. Tsai et al. (2019) proposed a crossmodal attention-based Transformer to provide tensor-level crossmodal adaptation that fuses multimodal information by directly attending to features in other modalities. Zadeh et al. (2019) developed a self-attention- and cross-attention-based Transformer to extract intra-modal and inter-modal emotional information, respectively. Li et al. (2022a) used crossmodal attention with a hierarchical structure to capture lexical features from different textual aspects for speech emotion recognition.

Despite these advances, some issues still remain challenging in multimodal fusion. First of all, dif-

ferent modalities may show mismatched affective tendencies, resulting in inter-modal incongruity – a general problem for MSER tasks. However, the majority of this topic is based on high-level comparison analysis between modalities, such as a person expressing praise while rolling his/her eyes (Wu et al., 2021). There is no evidence that such inter-modal incongruity can be tackled at the latent level by crossmodal attention. What’s more, to improve the performance of MSER tasks, certain modalities are usually fused repeatedly. Such an operation would bring information redundancy to the model and result in large model sizes, which hinder the real-world use of MSER. With these challenges in mind, we conduct a novel analysis on how cross-modal attention functions or fails, and propose a lightweight yet efficient model based on this.

## 3 Feasibility Analysis

Models utilizing data from different modalities usually outperform unimodal ones as more information is aggregated. Prior work has shown that learning with multiple modalities is superior to employing a subset of modalities, since the former has access to a better latent space representation (Huang et al., 2021). However, there is no guarantee that using multimodal data is always better than unimodal. For example, Huang et al. (2021) found that combining multiple modalities (text, audio, and video) underperforms the unimodal when sample sizes are relatively small. Moreover, Rajan et al. (2022) compared a self-attention and a crossmodal attention model for emotion recognition, showing no clear difference between the results of the two models.

As no evidence has been presented as to whether and why crossmodal attention works, we perform an analysis on the latent level to investigate how multimodal information interacts with each other and how inter-modal incongruity occurs. We conduct three experiments on CMU-MOSEI:

Exp 1. Investigates how source modality enhances target modality via crossmodal attention. We use the example of  $V \rightarrow T$  (text attended by vision). *Next, we hope to see how the combination of two modalities affects the third collectively.*

Exp 2. Investigates how the salient parts of the target modality are represented by self-attention with and without the combination of source modalities. We use the example of  $(A + V) \rightarrow T$  (text attended by cross-attention-fused audio-vision). *Further, we would like to know how different source*

<sup>1</sup>We will release code on paper acceptance.

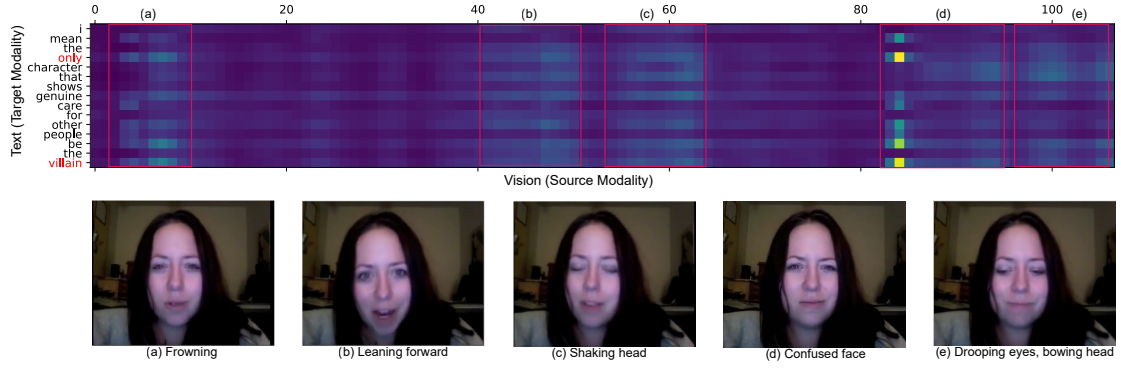


Figure 1: Exp 1. Heatmap of highlighted hidden states using crossmodal attention on vision (source) and text (target) modalities.

Crossmodal Attention: but this movie kind of came off like somebody really **sad** travel video, i **don't** know  
Text Self-attention: **but this movie kind of came off like somebody really sad travel video, i don't know**

Crossmodal Attention: and umm first of all, i, i actually was **never** a huge **fan** of the xx xxx tv series  
Text Self-attention: **and umm first of all, i, i actually was never a huge fan of the xx xxx tv series**

Crossmodal Attention: i really **love** the **melodramatic** over the top **crazy** action sort of thing going on  
Text Self-attention: **i really love the melodramatic over the top crazy action sort of thing going on**

Figure 2: Exp 2. Heatmap of highlighted words by using self-attention with and without crossmodal attention.

Crossmodal Attention ( $V \rightarrow T$ ): i **love** him as an **actor** based on his looks, as an actual **talented** actor, **no**  
Crossmodal Attention ( $A \rightarrow T$ ): i **love** him as an actor based on his looks, as an actual **talented** actor, **no**

Figure 3: Exp 3. Heatmap of highlighted words using different modalities ( $V$  or  $A$ ) in crossmodal attention.

modalities affect the target individually.

Exp 3. Investigates how the salient parts of the target modality are represented by crossmodal attention when using different source modalities. We use the examples of  $V \rightarrow T$  (text attended by vision) and  $A \rightarrow T$  (text attended by audio).

The experimental setup is shown in Table 1, and the visualization is shown in Figure 1, 2, and 3.

Table 1: Experimental setup for crossmodal attention analysis.  $T$  for Text,  $A$  for Audio, and  $V$  for Vision.

Exp.	Target modality	Source modality	Crossmodal	Self-attention
1	$T$	$V$	$V \rightarrow T$	$I$
2	$T$	$A + V$	$(A + V) \rightarrow T$	$T$
3	$T$	$A$ or $V$	$A \rightarrow T, V \rightarrow T$	$I$

Figure 1 shows the video frame (x-axis) and text words (y-axis). The salient emotional information captured by crossmodal attention is highlighted in the red box. It can be noticed that the highlighted parts are due to obvious facial or behavior changes of the character in the video, such as frowning or shaking head. The crossmodal attention successfully highlights the meaningful words associated with a facial expression (e.g., “only”, “villain”).

In Figure 2, it can be noted that when fused with

the combination of source modalities, text focuses more on the words related to emotional information with less noise from other words. For example, when with crossmodal attention, the word “sad” is the most salient in the first sentence, yet much less focused with self-attention. The same is true for the word “never” in the second sentence and the words “love” and “melodramatic” in the third sentence.

In Figure 3, we can see that when fused with different individual source modalities, the target modality (text) can be enhanced with disparate affective tendencies. When using vision as the source modality, the words “love” and “talented” are the most highlighted, representing a positive meaning. When using audio, however, “no” is the most focused word, showing negation is important. This phenomenon demonstrates that different modalities may contain mismatched affective tendencies. The existence of inter-modal incongruity has been found by high-level inter-modal comparison (Desai et al., 2022) and sentiment analysis (Li et al., 2019). Our finding demonstrates that the incongruity also exists at the latent level when using crossmodal attention, resulting in salient affective information

in one modality being distorted by the other.

Based on the above findings, we can find that crossmodal attention does help multimodal fusion by aligning two modalities to highlight the salient affective information in the target modality with complementary information from the source. According to the attention mechanism (Vaswani et al., 2017), this process can be described as mapping the Query (from the target) to the Key (from the source) and obtain scores for the Value (from the source). However, such a process could malfunction if the modalities have mismatched affective tendencies, which leaves the inter-modal incongruity difficult to resolve at the latent level.

#### 4 Proposed Approach – HCT-MG

To exploit the advantages of crossmodal attention while solving the above problems, we propose a new multimodal fusion approach: the Hierarchical Crossmodal Transformer with Modality Gating (HCT-MG), which is superior to existing methods in two aspects. First, while some previous studies treated all modalities equally and fused them in each step, leaving incongruity in the fusion (Tsai et al., 2019; Sahay et al., 2020), our HCT-MG fuses the auxiliary modalities first, leaving the primary for fusion in the final step. Second, while some previous studies determined a primary modality based on the hierarchy of modalities used (Rahman et al., 2020; Hazarika et al., 2020). Such a practice is empirical and due to the fixed hierarchy, the weighting pattern (e.g.,  $T \oplus W_1 A \oplus W_2 V$ ) cannot be changed during model training even though other modalities may become dominant. In contrast, HCT-MG automatically selects and dynamically changes the primary modality in each training batch and constructs the hierarchy accordingly, without worrying about which modalities are used. Thus, our proposed approach can remove incongruence and reduce redundancy while allowing the model to be modality agnostic.

The architecture is shown as Figure 4. HCT-MG is constructed based on three modalities: Text ( $T$ ), Audio ( $A$ ), and Vision ( $V$ ), and consists of four components: feature projection, modality gating, crossmodal Transformer, and weighted concatenation. Note that the modalities are not limited to  $T$ ,  $A$ , and  $V$ , as the modality gating enables to construct the best hierarchy for any three types of modality inputs.

**Feature Projection.** The input features are first

fed into 1D Convolutional (Conv1D) networks to integrate local contexts and project the features into the same hidden dimension. Then the features are passed to the Gated Recurrent Unit (GRU) networks, which encode global contexts by updating their hidden states recurrently and model the sequential structure of the extracted features accordingly. We use two sets of input features: one uses the same conventional feature extractors as (Tsai et al., 2018, 2019; Sahay et al., 2020) for fair comparison, while the other uses the same large pre-trained models as (Ando et al., 2023) for significant performance improvement. We will show the input features in the Appendix as they are not the focus of this work.

**Modality Gating.** MG determines which modality should be the primary one by its trainable weight for each modality during training, rather than by manual selection. Specifically, each modality is assigned a trainable weight whose value is based on its contribution to the final task, i.e., sentiment classification or emotion recognition. The larger the contribution of a modality, the larger its weight value. The sum of all trainable weights ( $\omega_p$ ,  $\omega_{a_1}$ , and  $\omega_{a_2}$ ) equals to 1, and we allow the weights to be updated in every training batch to ensure that modality gating can be well adapted to any type of input modality. We will discuss how modality gating works in Sec. 5.3.2.

**Crossmodal Transformer.** As a variant of self-attention, cross-attention (Lu et al., 2019) transforms the signals from the source modality into a different set of Key-Value pairs to interact with the target modality. This has proven useful in various domains (Zhang et al., 2022; Rashed et al., 2022). The crossmodal Transformer used here is the same as MulT (Tsai et al., 2019), which is a deep stacking of several crossmodal attention blocks with layer normalization and positional embeddings. Unlike MulT, which has six crossmodal Transformers implemented in the same step, we use two crossmodal Transformers in the first step to obtain enhanced auxiliary modalities:

$$\hat{A}_1 = CMT(A_2 \rightarrow A_1) \quad (1)$$

$$\hat{A}_2 = CMT(A_1 \rightarrow A_2) \quad (2)$$

Then in the second step, another two crossmodal Transformers are used to yield the enhanced primary modality representations:

$$\hat{P}_{\hat{A}_1} = CMT(\hat{A}_1 \rightarrow P) \quad (3)$$

$$\hat{P}_{\hat{A}_2} = CMT(\hat{A}_2 \rightarrow P) \quad (4)$$



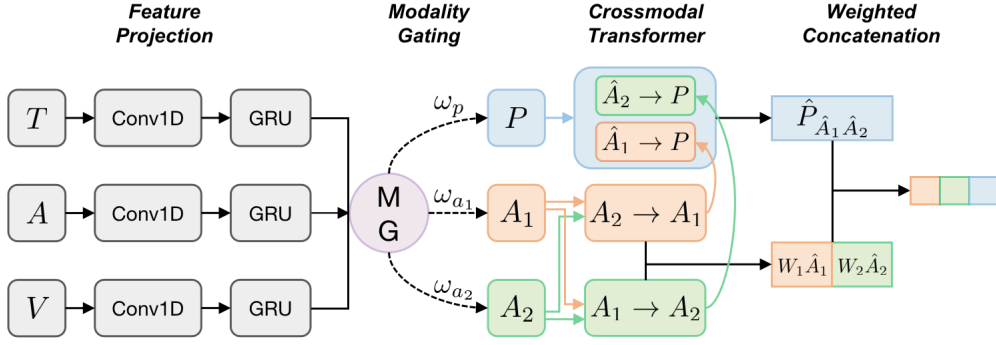


Figure 4: Architecture of HCT-MG.

**Weighted Concatenation.** After obtaining the enhanced  $\hat{P}_{\hat{A}_1}$  and  $\hat{P}_{\hat{A}_2}$ , we concatenate them and use the self-attention to find its salient parts as the final primary representation:

$$\hat{P}_{\hat{A}_1\hat{A}_2} = SA(Concat[\hat{P}_{\hat{A}_1}; \hat{P}_{\hat{A}_2}]) \quad (5)$$

By now, crossmodal representations of every modality have been generated:  $\hat{A}_1$ ,  $\hat{A}_2$ , and  $\hat{P}_{\hat{A}_1\hat{A}_2}$ . We concatenate them for the final representation:

$$Z = Concat[W_1\hat{A}_1; W_2\hat{A}_2; \hat{P}_{\hat{A}_1\hat{A}_2}] \quad (6)$$

where  $W_1$  and  $W_2$  are the weight matrices, which are learned by the model itself to control how much information to extract from the two auxiliary modalities.

## 5 Experiments

We describe the datasets and report our results via a comparison with prior models. Because MG will become obsolete once the primary modality selection converges, we freeze it at that point. We present the process by which HCT-MG selects the primary modality in the ablation study.

### 5.1 Datasets and Evaluation Metrics

**CMU-MOSI** and **CMU-MOSEI** are sentiment analysis datasets containing video clips from YouTube, annotated with sentiment scores in the range of  $[-3, 3]$ . The former has 2,199 samples, while the latter has 23,454. **IEMOCAP** is a multi-modal dataset for emotion recognition. Following prior work, we use four emotions (happy, sad, angry, and neutral) for the experimental evaluation, bringing in 4,453 samples.

As with prior work on MOSI and MOSEI, we evaluate the performances using the following metrics: 7-class accuracy (i.e., Acc-7: sentiment score classification in the same scale as the labeled scores); binary accuracy (i.e., Acc-2: positive/negative sentiment polarity); F1 score; Mean

Absolute Error (MAE); and the correlation of the recognition results with ground truth. On IEMOCAP, we report the binary classification accuracy (one versus the others) and F1 score.

### 5.2 Experimental Evaluation

We use the features provided by the CMU-SDK (Zadeh et al., 2018c), which splits the datasets into folds for training, validation, and testing. The experimental settings are presented in the Appendix.

#### 5.2.1 Baselines

We perform a comparative study against our approach, considering two aspects: 1) mainstream models that have been widely compared; 2) lightweight models with similar sizes to ours. Note that not every baseline has been evaluated on all three datasets. The baselines are as below:

Early Fusion LSTM (**EF-LSTM**) and Late Fusion LSTM (**LF-LSTM**) (Tsai et al., 2018). Attention or Transformer-based fusion: **RAVEN** (Wang et al., 2019), **MuT** (Tsai et al., 2019). Graph-based fusion: **Graph-MFN** (Zadeh et al., 2018b). Low-rank-based fusion: **LMF** (Liu et al., 2018). Cyclic translations-based fusion: **MCTN** (Pham et al., 2019). Context-aware attention-based fusion: **CIA** (Chauhan et al., 2019). Multi-attention Recurrent-based fusion: **MARN** (Zadeh et al., 2018c). Temporal memory-based fusion: **MFN** (Zadeh et al., 2018a). Recurrent multiple stages-based fusion: **RMFN** (Liang et al., 2018). Low-rank Transformer-based fusion: **LMF-MuT** (Sahay et al., 2020). We also include several of the above-mentioned models enhanced by Connectionist Temporal Classification (CTC) (cf. Tsai et al. (2019)). As the recognition of fine-grained emotions is significantly improved by word alignment, we do not include EF-LSTM, RAVEN, MCTN (which are tested on aligned corpora) for comparisons on IEMOCAP. Because our crossmodal

Transformer component is based on MulT, we reproduced it in the same experimental environment for a fair comparison.

## 5.2.2 Results

Table 2: Comparison results on CMU-MOSI, MOSEI, and IEMOCAP for sentiment analysis and emotion recognition.  $\uparrow$ : higher is better;  $\downarrow$ : lower is better; \*: reproduced from the official code.

Models	CMU-MOSI				
	Acc-7 $\uparrow$	Acc-2 $\uparrow$	F1-score $\uparrow$	Corr $\uparrow$	MAE $\downarrow$
EF-LSTM	33.7	75.3	75.2	0.608	1.023
RAVEN	33.2	78.0	76.6	0.691	0.915
MCTN	35.6	79.3	79.1	0.676	0.909
CTC+EF-LSTM	31.0	73.6	74.5	0.542	1.078
CTC+RAVEN	31.7	72.7	73.1	0.544	1.076
CTC+MCTN	32.7	75.9	76.4	0.613	0.991
MARN	34.7	77.1	77.0	0.625	0.968
MFN	34.1	77.4	77.3	0.632	0.965
RMFN	38.3	78.4	78.0	0.681	0.922
LMF	32.8	76.4	75.7	0.668	0.912
CIA	38.9	79.8	79.5	0.689	0.914
LF-LSTM(1.24M)	33.7	77.6	77.8	0.624	0.988
MulT*(1.07M)	34.3	80.3	80.4	0.645	1.008
LMF-MulT(0.84M)	34.0	78.5	78.5	0.681	0.957
HCT-MG(0.54M)	<b>39.4</b>	<b>82.5</b>	<b>82.5</b>	<b>0.710</b>	<b>0.881</b>

Models	CMU-MOSEI				
	Acc-7 $\uparrow$	Acc-2 $\uparrow$	F1-score $\uparrow$	Corr $\uparrow$	MAE $\downarrow$
EF-LSTM	47.4	78.2	77.9	0.642	0.616
RAVEN	50.0	79.1	79.5	0.662	0.614
MCTN	49.6	79.8	80.6	0.670	0.609
CTC+EF-LSTM	46.3	76.1	75.9	0.585	0.680
CTC+RAVEN	45.5	75.4	75.7	0.599	0.664
CTC+MCTN	48.2	79.3	79.7	0.645	0.631
LMF	48.0	<b>82.0</b>	<b>82.1</b>	0.677	0.623
Graph-MFN	45.0	76.9	77.0	0.540	0.710
CIA	50.1	80.4	78.2	0.590	0.680
LF-LSTM(1.24M)	48.8	77.5	78.2	0.656	0.624
MulT*(1.07M)	50.4	80.7	80.6	0.677	0.617
LMF-MulT(0.84M)	49.3	80.8	81.3	0.668	0.620
HCT-MG(0.78M)	<b>50.6</b>	81.8	81.9	<b>0.691</b>	<b>0.593</b>

Models	IEMOCAP							
	Happy		Sad		Angry		Neutral	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
CTC+EF-LSTM	76.2	75.7	70.2	70.5	72.7	67.1	58.1	57.4
CTC+RAVEN	77.0	76.8	67.6	65.6	65.0	64.1	<b>62.0</b>	<b>59.5</b>
CTC+MCTN	80.5	77.5	72.0	<b>71.7</b>	64.9	65.6	49.4	49.3
LF-LSTM(1.24M)	72.5	71.8	72.9	70.4	68.6	<b>67.9</b>	59.6	56.2
MulT*(1.07M)	<b>85.6</b>	<b>79.0</b>	<b>79.4</b>	70.3	<b>75.8</b>	65.4	59.5	44.7
LMF-MulT(0.86M)	<b>85.6</b>	<b>79.0</b>	<b>79.4</b>	70.3	<b>75.8</b>	65.4	59.2	44.0
HCT-MG(0.55M)	<b>85.6</b>	<b>79.0</b>	<b>79.4</b>	70.3	<b>75.8</b>	65.4	61.0	50.5

The comparison results are shown in Table 2. On MOSI and MOSEI, it can be seen that HCT-MG improves on most metrics compared to the baselines, no matter the aligned or unaligned settings. Besides, HCT-MG outperforms the models that have similar sizes, particularly on MOSI. On IEMOCAP, we achieve the overall best performance with the same or similar results on most metrics on *Happy*, *Sad*, and *Angry*. On *Neutral*, CTC+RAVEN outperforms the others on F1 but is much worse on the other three emotions, which means it is better in distinguishing emotion from non-emotion

than recognizing fine-grained emotions. Moreover, compared to LMF-MulT and MulT, which used the same crossmodal Transformer as ours, HCT-MG has the same results on *Happy*, *Sad*, and *Angry*. It is likely due to the performance bottleneck using this crossmodal Transformer without larger model sizes. However, HCT-MG outperforms them on *Neutral*, which is plausible as our hierarchical architecture works well on distinguishing all emotions.

Moreover, HCT-MG was used as the crossmodal attention  $(A + V) \rightarrow T$  in Exp. 2 (in Sec. 3). Figure 2 clearly shows that our approach strengthens the most salient affective parts effectively. Furthermore, we present some examples where incongruity or ambiguity exists. In Table 3, MulT fails to handle these difficult cases (videos available<sup>2</sup>), producing results that contradict ground truths. In contrast, our approach can recognize true sentiments with very close scores. The examples demonstrate that HCT-MG can successfully integrate auxiliary modalities with the primary one. Their heatmaps are presented in the Appendix.

## 5.3 Further Analysis and Discussion

To verify that our approach alleviates the inter-modal incongruity issue and to demonstrate how MG dynamically changes the primary modality to construct according hierarchy, we conducted the following studies.

### 5.3.1 Resolution of Inter-Modal Incongruity

As shown in Figure 3, the affective words are enhanced by the auxiliary modalities. However, the vision modality focuses more on positive words while the audio highlights negation the most, which likely changes affective tendency. We reimplemented MulT and extracted the attention of its enhanced text modality (attended by audio and vision) to compare with ours. In Figure 5, it can be seen that the positive and negative words are treated equally by MulT, which leaves the inter-modal incongruity unsolved. This is because MulT fuses audio and vision with text at the same level and simply concatenates two enhanced text modalities. On the other hand, our approach barely gives attention to the word “no”, showing that the incongruity is resolved at the latent level (Yet the sentiment score was labeled as *-1*. See Limitations). The examples in Table 3 also demonstrate that the incongruity and ambiguity are largely resolved by our approach.

<sup>2</sup><https://sites.google.com/view/acsubmission>

Table 3: Examples containing incongruity or ambiguity from CMU-MOSI.

#	Spoken words + acoustic and visual behaviors	Ground truth	MuT	Ours
1	"And that's why I was not excited about the fourth one." + Uninterested tone and facial expression	-1.4	1.185	-1.416
2	"I give Shrek Forever After directed by Mike Mitchell a grade of B minus." + Smile face	1.0	-0.576	0.959
3	"Um in general um, the little kids seemed to like it that were in there." + Skeptical tone and facial expression	0.8	-1.151	0.700
4	"I honestly want the aliens to win." + Negative tone (somewhat disdainful)	-1.6	0.995	-1.906

MuT: i love him as an actor based on his looks, as an actual talented actor, no  
Ours: i love him as an actor based on his looks, as an actual talented actor, no

Figure 5: A comparison example of incongruity.

### 5.3.2 Automatic Modality Selection by MG

As the MG automatically selects the primary modality by adjusting the weights for each modality, we show how the weights vary during training using MOSI. The weight of a modality denotes the probability that this modality is selected as the primary one. Figure 6 shows how the probabilities of the modalities vary in the first epoch. It can be noted that the text modality is not the primary one at the beginning but gradually dominates after batch 60 during the training. Figure 7 shows the variation in the average probability of each modality with epoch during training. It can be seen that the text modality does dominate, and the probability distribution starts to converge at around 40 epochs. The results of Figure 6 and 7 are consistent with Table 4. It confirms again that our empirical approach is well-founded. (Note that, the three modalities had very close probabilities on IEMOCAP (slightly above 0.33), which is in line with their results in Table 4. It is likely due to the fact that IEMOCAP was collected in a lab environment where emotions can be well expressed in all modalities.) Furthermore, now that we know how HCT-MG works, we expects its utility expand beyond more general scenarios with different and more signals (e.g., Electroencephalogram (EEG)).

### 5.3.3 Ablation Study on the Removal of MG

Although MG automatically selects text as the primary modality, we would like to see if this phenomenon really brings the best performance compared to prior work whose hierarchies are fixed. Thus, we removed MG and constructed an HCT model with three hierarchies, where  $T$ ,  $A$ , and  $V$  were manually selected as the primary modality, respectively, for performance comparison. As shown in Table 4, selecting text as the primary modality

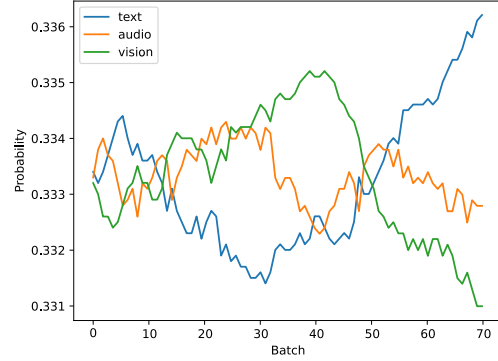


Figure 6: Probability (weight) variation of each modality in the first epoch.

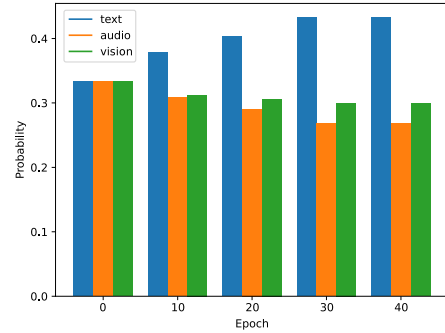


Figure 7: Average probability (weight) variation of each modality with epoch.

achieves the best performance on every metric on MOSI and MOSEI and on most metrics on IEMOCAP, verifying the rationality and efficacy of the automatic modality selection by MG.

Moreover, although the results of selecting  $T$  as the primary modality are the best among the three, the scores are still lower than those of HCT-MG in Table 2. This phenomenon is reasonable because the primary modality kept changing during training, even though  $T$  is primary overall. Such a dynamic property is encoded by the HCT-MG and thus yields better results.

### 5.3.4 Discussion on Modality Selection

To the best of our knowledge, there is no detailed explanation as to why choosing  $T$  as the primary modality works best for multimodal language anal-

Table 4: Performance comparison by manually selecting different primary modalities.  $\uparrow$ : higher is better;  $\downarrow$ : lower is better.

Primary modality	CMU-MOSI				
	Acc-7 $\uparrow$	Acc-2 $\uparrow$	F1 $\uparrow$	Corr $\uparrow$	MAE $\downarrow$
<i>T</i>	<b>38.9</b>	<b>82.5</b>	<b>82.6</b>	<b>0.717</b>	<b>0.859</b>
<i>A</i>	37.5	81.3	81.3	0.705	0.883
<i>V</i>	38.3	80.9	81.0	0.679	0.909

Primary modality	CMU-MOSEI				
	Acc-7 $\uparrow$	Acc-2 $\uparrow$	F1 $\uparrow$	Corr $\uparrow$	MAE $\downarrow$
<i>T</i>	<b>50.1</b>	<b>81.8</b>	<b>81.9</b>	<b>0.685</b>	<b>0.601</b>
<i>A</i>	47.5	79.6	80.3	0.650	0.644
<i>V</i>	48.7	80.8	81.0	0.659	0.633

Primary modality	IEMOCAP							
	Happy		Sad		Angry		Neutral	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
<i>T</i>	<b>85.6</b>	<b>79.4</b>	<b>79.5</b>	<b>70.6</b>	75.8	65.4	59.6	<b>55.6</b>
<i>A</i>	85.4	<b>79.4</b>	78.8	70.4	75.4	65.5	59.3	52.4
<i>V</i>	<b>85.6</b>	79.2	79.4	70.3	<b>75.9</b>	<b>65.6</b>	<b>60.4</b>	53.3

ysis, especially MSER. Here, we gather the following empirical findings from different perspectives.

1) There is a clear temporal pattern when people express emotions via vision and audio modalities: visual signals usually precede audio by around 120ms (Grant and Greenberg, 2001). 2) People can behave quite differently from what they say in spoken dialogues. For example, positive behaviors sometimes come along with a negative sentence to ease the embarrassment (Li et al., 2019), and a positive sentence can be said in a negative way to express sarcasm (Castro et al., 2019). 3) In MSER applications, a misrecognition of emotion at the level of sentiment polarity would lead to a fatal error (imagine that the system responds "Good to hear that!" in a happy voice when the user in fact feels sad). On the other hand, misclassifying an emotion as another that has the same polarity may well be tolerable (Tokuhisa et al., 2008; Li, 2018). 4) Sentiment largely depends on textual information (Lindquist et al., 2015). Using text as the primary modality in fine-tuning and shifting the language-only position of a word to the new position in light of audio-visual information allows the language models (e.g., BERT, XLNet) to better yield sentiment scores (Rahman et al., 2020). 5) Modality refers to the way in which something is expressed or perceived. Unlike audio and vision, which are raw (low-level) modalities closest to sensors, text is a relatively abstract and high-level modality that is farther from sensors (Baltrušaitis et al., 2018). According to the nature of the human brain's hierarchical perceptual processing, low-level information is processed first, followed by high-level

information (Peelle et al., 2010). Thus, using a hierarchical model to process low-level features and fuse high-level ones sequentially can yield better representations for emotion recognition (Tian et al., 2016; Li et al., 2022a).

We expect that our proposed HCT-MG will bring new insights to the literature and, together with the aforementioned studies, provide a theoretical basis to support that *text is relatively independent from audio and vision but significantly contributes to affective polarity*.

## 6 Conclusions

In this work, we analyzed crossmodal attention-based multimodal fusion and propose a hierarchical crossmodal Transformer with modality gating for incongruity-aware multimodal sentiment and emotion recognition. The major contributions are:

1) We demonstrate the existence of inter-modal incongruity at the latent level due to crossmodal attention. Specifically, we show that crossmodal attention can help to capture affective information across modalities and enhance salient parts in the target modality, but it can also bring mismatched affective tendencies from different modalities.

2) We propose a hierarchical crossmodal Transformer with modality gating – HCT-MG, that automatically selects the primary modality during training. This model requires fewer fusions and does not repeatedly fuse a single modality, reducing the model size to less than 1M while significantly outperforming existing models of similar size.

3) We further analyze the mechanism and feasibility of automatic modality selection by MG and show that the selection process supports the primacy of text in previous sentiment and emotion analysis studies, adding new insights to the literature.

In our future work, we will test HCT-MG in other domains, especially where the affective tendencies of different modalities easily mismatch, such as humor and sarcasm detection (Pramanick et al., 2022), and where other modalities exist, such as biological signals (Li et al., 2022b). We will also try dimensionality reduction techniques (Shao et al., 2022) to further reduce redundant and misleading information from learned representations. We expect HCT-MG could be used for solving multimodal affective computing tasks in real life.



## Limitations

Although the inter-modal incongruity is largely removed by the hierarchical architecture, the affective tendency could be wrong if ambiguity exists in the primary modality. In Figure 5, our approach recognizes this sample as positive with the score of  $+1.84$ , but the ground truth is labeled as negative with the score of  $-1$ . It is likely because the first half of the text denotes a positive sentiment, yet the second is obviously negative. Without contextual knowledge, it is almost impossible for a system to know that the second half is the focus of the content, as humans can. As our model uses GloVe as text features for a fair comparison with prior work, which are non-contextual word embeddings, it is hard to capture the semantic meaning. This problem can likely be resolved with contextual language models (e.g., BERT) or auxiliary modalities that have the same affective tendency (which can modify the primary modality).

Also, as with all the other supervised learning tasks, our model relies on the accuracy of the labels. However, it is not easy to label difficult cases where inter-modal incongruity, ambiguous emotions, or missing information exist. Take the cases in Table 3 as an example: 1) In our opinion, #2 should be a neutral emotion with the value of 0, as the person is just stating a fact. 2) There is a word “but” missing at the end of the labeled sentence of #3 (can be clearly noticed in the audio or the video), which is a sign indicating a turnaround in attitude. It can be regarded as the same situation as the sample in Figure 5. However, #3 separates the whole sentence into two sub-sentences, while the sample in Figure 5 combines two sub-sentences as a whole. Such an inconsistency in labeling introduces incongruity and ambiguity into the tasks and hinders the training of robust and applicable models.

## References

- Atsushi Ando, Ryo Masumura, Akihiko Takashima, Satoshi Suzuki, Naoki Makishima, Keita Suzuki, Takafumi Moriya, Takanori Ashihara, and Hiroshi Sato. 2023. On the use of modality-specific large-scale pre-trained encoders for multimodal sentiment analysis. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 739–746. IEEE.
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443.
- Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335–359.
- Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. 2019. Towards multimodal sarcasm detection (an \_obviously\_ perfect paper). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4619–4629.
- Dushyant Singh Chauhan, Md Shad Akhtar, Asif Ekbal, and Pushpak Bhattacharyya. 2019. Context-aware interactive attention for multi-modal sentiment and emotion analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5647–5657.
- Ze-Jing Chuang and Chung-Hsien Wu. 2004. Multimodal emotion recognition from speech and text. In *International Journal of Computational Linguistics & Chinese Language Processing, Volume 9, Number 2, August 2004: Special Issue on New Trends of Speech and Language Processing*, pages 45–62.
- Gilles Degottex, John Kane, Thomas Drugman, Tuomo Raitio, and Stefan Scherer. 2014. Covarep—a collaborative voice analysis repository for speech technologies. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 960–964. IEEE.
- Poorav Desai, Tanmoy Chakraborty, and Md Shad Akhtar. 2022. Nice perfume. how long did you marinate in it? multimodal sarcasm explanation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10563–10571.
- Ken W Grant and Steven Greenberg. 2001. Speech intelligibility derived from asynchronous processing of auditory-visual information. In *AVSP 2001-International Conference on Auditory-Visual Speech Processing*.
- Devamanyu Hazarika, Roger Zimmermann, and Soujanya Poria. 2020. Misa: Modality-invariant and-specific representations for multimodal sentiment analysis. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1122–1131.
- Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. 2021. What makes multi-modal learning better than single (provably). *Advances in Neural Information Processing Systems*, 34:10944–10956.
- Yuanchao Li. 2018. Towards improving speech emotion recognition for in-vehicle agents: Preliminary results of incorporating sentiment analysis by using early and late fusion methods. In *Proceedings of*

706	<i>the 6th International Conference on Human-Agent Interaction</i> , pages 365–367.	760
707		761
708	Yuanchao Li, Peter Bell, and Catherine Lai. 2022a. Fus-	762
709	ing ASR outputs in joint training for speech emotion	763
710	recognition. In <i>ICASSP 2022-2022 IEEE Interna-</i>	764
711	<i>tional Conference on Acoustics, Speech and Signal</i>	
712	<i>Processing (ICASSP)</i> , pages 7362–7366. IEEE.	
713	Yuanchao Li, Peter Bell, and Catherine Lai. 2022b.	765
714	Multimodal dyadic impression recognition via lis-	766
715	tener adaptive cross-domain fusion. <i>arXiv preprint</i>	767
716	<i>arXiv:2211.05163</i> .	768
717	Yuanchao Li, Carlos Toshinori Ishi, Koji Inoue, Shizuka	769
718	Nakamura, and Tatsuya Kawahara. 2019. Express-	770
719	ing reactive emotion based on multimodal emotion	771
720	recognition for natural conversation in human–robot	
721	interaction. <i>Advanced Robotics</i> , 33(20):1030–1041.	
722	Yuanchao Li, Tianyu Zhao, and Xun Shen. 2020.	772
723	Attention-based multimodal fusion for estimating hu-	773
724	man emotion in real-world HRI. In <i>Companion of</i>	774
725	<i>the 2020 ACM/IEEE International Conference on</i>	775
726	<i>Human-Robot Interaction</i> , pages 340–342.	776
727	Paul Pu Liang, Ziyin Liu, Amir Zadeh, and Louis-	777
728	Philippe Morency. 2018. Multimodal language anal-	
729	ysis with recurrent multistage fusion. In <i>Proceedings</i>	
730	<i>of the 2018 Conference on Empirical Methods in</i>	
731	<i>Natural Language Processing</i> , pages 150–161.	
732	Kristen A Lindquist, Jennifer K MacCormack, and	778
733	Holly Shablack. 2015. The role of language in emo-	779
734	tion: Predictions from psychological constructionism.	780
735	<i>Frontiers in psychology</i> , 6:444.	781
736	Zhun Liu, Ying Shen, Varun Bharadhwaj Lakshmi-	782
737	narasimhan, Paul Pu Liang, Amir Zadeh, and Louis-	
738	Philippe Morency. 2018. Efficient low-rank multi-	
739	modal fusion with modality-specific factors. <i>arXiv</i>	
740	<i>preprint arXiv:1806.00064</i> .	
741	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee.	783
742	2019. Vilbert: Pretraining task-agnostic visiolinguis-	784
743	tic representations for vision-and-language tasks. <i>Ad-</i>	785
744	<i>vances in neural information processing systems</i> , 32.	786
745	Jonathan E Peelle, Ingrid Johnsrude, and Matthew H	
746	Davis. 2010. Hierarchical processing for speech in	
747	human auditory cortex and beyond. <i>Frontiers in hu-</i>	
748	<i>man neuroscience</i> , page 51.	
749	Jeffrey Pennington, Richard Socher, and Christopher D	787
750	Manning. 2014. Glove: Global vectors for word rep-	788
751	resentation. In <i>Proceedings of the 2014 conference</i>	789
752	<i>on empirical methods in natural language processing</i>	790
753	<i>(EMNLP)</i> , pages 1532–1543.	
754	Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-	791
755	Philippe Morency, and Barnabás Póczos. 2019.	792
756	Found in translation: Learning robust joint repre-	793
757	sentations by cyclic translations between modalities.	794
758	In <i>Proceedings of the AAAI Conference on Artificial</i>	795
759	<i>Intelligence</i> , volume 33, pages 6892–6899.	
	Shraman Pramanick, Aniket Roy, and Vishal M Patel.	796
	2022. Multimodal learning using optimal transport	797
	for sarcasm and humor detection. In <i>Proceedings of</i>	798
	<i>the IEEE/CVF Winter Conference on Applications of</i>	799
	<i>Computer Vision</i> , pages 3930–3940.	800
	Wasifur Rahman, Md Kamrul Hasan, Sangwu Lee,	
	Amir Zadeh, Chengfeng Mao, Louis-Philippe	
	Morency, and Ehsan Hoque. 2020. Integrating multi-	
	modal information in large pretrained transformers.	
	In <i>Proceedings of the conference. Association for</i>	
	<i>Computational Linguistics. Meeting</i> , volume 2020,	
	page 2359. NIH Public Access.	
	Vandana Rajan, Alessio Brutti, and Andrea Cavallaro.	801
	2022. Is cross-attention preferable to self-attention	802
	for multi-modal emotion recognition? In <i>ICASSP</i>	803
	<i>2022-2022 IEEE International Conference on Acous-</i>	804
	<i>tics, Speech and Signal Processing (ICASSP)</i> , pages	805
	4693–4697. IEEE.	806
	Ahmed Rashed, Shereen Elsayed, and Lars Schmidt-	807
	Thieme. 2022. Context and attribute-aware sequen-	
	tial recommendation via cross-attention. In <i>Proceed-</i>	
	<i>ings of the 16th ACM Conference on Recommender</i>	
	<i>Systems</i> , pages 71–80.	
	Saurav Sahay, Eda Okur, Shachi H Kumar, and Lama	808
	Nachman. 2020. Low rank fusion based trans-	809
	formers for multimodal sequences. <i>arXiv preprint</i>	810
	<i>arXiv:2007.02038</i> .	811
	Shun Shao, Yftah Ziser, and Shay B Cohen. 2022. Gold	
	doesn’t always glitter: Spectral removal of linear	
	and nonlinear guarded attribute information. <i>arXiv</i>	
	<i>preprint arXiv:2203.07893</i> .	
	Leimin Tian, Johanna Moore, and Catherine Lai. 2016.	
	Recognizing emotions in spoken dialogue with hi-	
	erarchically fused acoustic and lexical features. In	
	<i>2016 IEEE Spoken Language Technology Workshop</i>	
	<i>(SLT)</i> , pages 565–572. IEEE.	
	Ryoko Tokuhisa, Kentaro Inui, and Yuji Matsumoto.	812
	2008. Emotion classification using massive examples	813
	extracted from the web. In <i>Proceedings of the 22nd</i>	814
	<i>International Conference on Computational Linguis-</i>	815
	<i>tics (Coling 2008)</i> , pages 881–888.	816
	Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang,	
	J Zico Kolter, Louis-Philippe Morency, and Ruslan	
	Salakhutdinov. 2019. Multimodal transformer for	
	unaligned multimodal language sequences. In <i>Pro-</i>	
	<i>ceedings of the conference. Association for Computa-</i>	
	<i>tional Linguistics. Meeting</i> , volume 2019, page 6558.	
	NIH Public Access.	
	Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh,	
	Louis-Philippe Morency, and Ruslan Salakhutdinov.	
	2018. Learning factorized multimodal representa-	
	tions. <i>arXiv preprint arXiv:1806.06176</i> .	
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	
	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	
	Kaiser, and Illia Polosukhin. 2017. Attention is all	
	you need. <i>Advances in neural information processing</i>	
	<i>systems</i> , 30.	

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*, volume 33, pages 7216–7223.

Yang Wu, Yanyan Zhao, Xin Lu, Bing Qin, Yin Wu, Jian Sheng, and Jinlong Li. 2021. Modeling incongruity between modalities for multimodal sarcasm detection. *IEEE MultiMedia*, 28(2):86–95.

Nan Xu, Wenji Mao, and Guandan Chen. 2018. A co-memory network for multimodal sentiment analysis. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 929–932.

Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. Tensor fusion network for multimodal sentiment analysis. *arXiv preprint arXiv:1707.07250*.

Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018a. Memory fusion network for multi-view sequential learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Amir Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018b. Multi-modal 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.

Amir Zadeh, Paul Pu Liang, Soujanya Poria, Prateek Vij, Erik Cambria, and Louis-Philippe Morency. 2018c. Multi-attention recurrent network for human communication comprehension. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

Amir Zadeh, Chengfeng Mao, Kelly Shi, Yiwei Zhang, Paul Pu Liang, Soujanya Poria, and Louis-Philippe Morency. 2019. Factorized multimodal transformer for multimodal sequential learning. *arXiv preprint arXiv:1911.09826*.

Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. 2016. Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages. *IEEE Intelligent Systems*, 31(6):82–88.

Shucong Zhang, Malcolm Chadwick, Alberto Gil CP Ramos, and Sourav Bhattacharya. 2022. Cross-attention is all you need: Real-time streaming transformers for personalised speech enhancement. *arXiv preprint arXiv:2211.04346*.

## A Appendix

### A.1 Datasets

Table 5: Data distribution and modality sampling rate of CMU-MOSI and CMU-MOSEI.  $S_A$  for audio sampling rate and  $S_V$  for vision sampling rate.

Dataset	Train	Valid	Test	Total	$S_A$	$S_V$
CMU-MOSI	1284	229	686	2199	12.5	15
CMU-MOSEI	16,326	1871	4659	22,856	20	15

Table 6: Data distribution of four emotions in the IEMO-CAP dataset.

Emotions	Train	Valid	Test	Total
Neural	954	358	383	1695
Happy	338	116	135	589
Sad	690	188	193	1071
Angry	735	136	227	1098
Total	2717	798	938	4453

### A.2 Extracted Features

The sequence lengths and feature dimensions of the three modalities in the three benchmarks are shown in Table 7.

Table 7: Sequence lengths and feature dimensions of the three modalities in the three benchmark datasets. \*: The development team screened the vision and audio features of CMU-MOSI.

Dataset	Text		Vision		Audio	
	len	dim	len	dim	len	dim
CMU-MOSI	50	300	500	20*	375	5*
CMU-MOSEI	50	300	500	35	500	74
IEMOCAP	20	300	500	35	400	74

**Textual Features: GloVe.** The transcriptions in all three datasets use the global word embeddings generated by GloVe. This distributed representation allows words in the same context to be close to each other in the vector space and maintain specific relationships (Pennington et al., 2014). For this pre-extracted data, the text modal features are trained and derived from 840 billion tokens with 300 dimensions of GloVe embeddings.

**Vision Features: FACET.** FACET is a commercial facial emotion detection software developed by iMotions<sup>3</sup>. The software can demonstrate 35 facial action units and record facial muscle movements to represent frame-by-frame emotions.

**Audio Features: COVAREP.** COVAREP is an open-source repository for speech processing, supporting collaboration and free access. The features

<sup>3</sup><https://imotions.com/platform/>

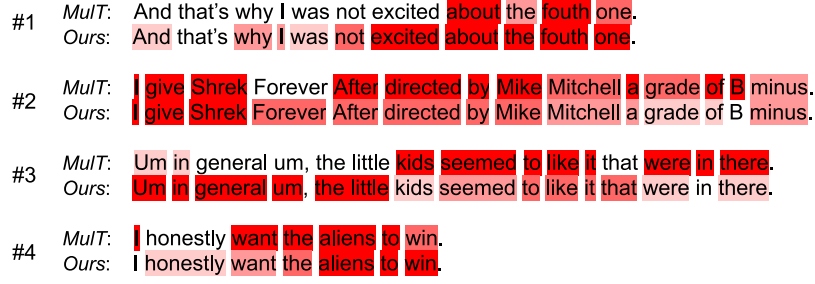


Figure 8: Heatmap comparison of the examples in Table 3 (Page 7).

of the processed speech data are based on pitch tracking, polarity detection, spectral envelopes, glottal flow, and other common speech features (Degottex et al., 2014). The pre-extracted data contains 74 dimensions of speech features.

### A.3 Hyperparameters Tuning

After tuning the hyperparameters, we find the optimal settings, as shown in Table 8.

Table 8: Hyperparameter settings for the three datasets.

Setting	CMU-MOSI	CMU-MOSEI	IEMOCAP
learning rate	1e-3	1e-3	1e-5
batch size	36	64	16
hidden size	40	40	40
kernel (T/A/V)	1/1/1	1/1/1	1/1/1
decay when	20	20	20
number of epochs	30	30	60
transformer layers	2	4	2
attention heads	5	5	5

### A.4 Heatmap Comparison

The heatmap values are from the enhanced text modalities (attended by audio and vision) of MulT and ours. With the the hierarchical architecture of our approach, some words that are not highlighted in MulT are highlighted.