

Customized-Allocation MoE: Reliable Dynamic Combination via the Attraction-Repulsion for Multimodal Sentiment Analysis

Anonymous ACL submission

Abstract

Multimodal Sentiment Analysis (MSA) aims to infer sentiment by jointly leveraging textual, visual, and acoustic modalities. However, a core challenge remains: *How to dynamically identify and leverage sample-level dependency preferences among different modality combinations?* To address this, we decompose the challenge into two subproblems: combination matching and fusion validation. Correspondingly, we propose the Customized-Allocation Mixture-of-Experts (CA-MoE), a novel framework that consists of two core complementary components that enable dynamic and sample-level modality routing within the MoE architecture. First, Affinity-guided Customized Modality Allocation (ACMA) acts as a distributor and leverages Geometry-Gradient Affinity (G^2 -Affinity) to guide an attraction-repulsion routing mechanism for customized allocation of modality combinations. After expert fusion, we then design Reliability-Aware Expert Selection (RAES) to jointly consider sentiment angular-prototype proximity and competitive magnitude intensity of representation. This yields a reliability selection matrix that weights over experts for the final sentiment prediction. Extensive experiments on three benchmark MSA datasets demonstrate that CA-MoE achieves significant or competitive performance over state-of-the-art methods.

1 Introduction

Human emotional expression is inherently multimodal, interweaving spoken language, facial expressions, and vocal cues. Therefore, unimodal approaches inevitably lose complementary affective signals, leading to incomplete sentiment understanding. (Baltruaitis et al., 2019). Multimodal Sentiment Analysis (MSA) aims to predict sentiment by jointly modeling text, vision, and audio, to reduce modality-specific bias that arises when any single modality is used alone. (Poria et al.,

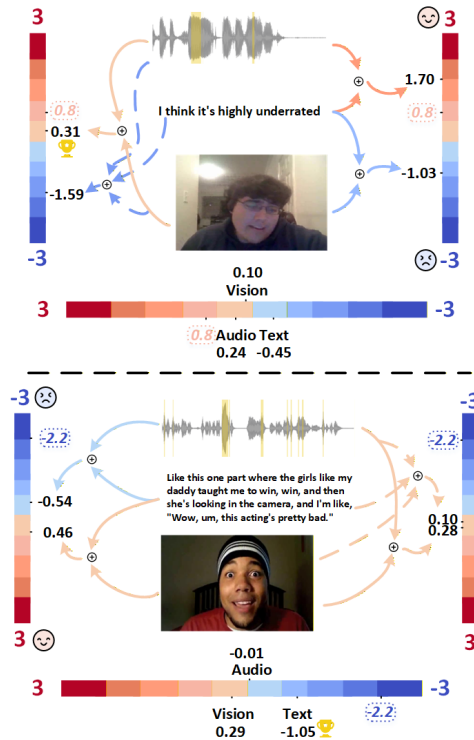


Figure 1: Illustrative example of two samples in dataset of CMU-MOSI (Zadeh et al., 2016) with widely-recognized MISA (Hazarika et al., 2020), which quantitatively demonstrates that predictions obtained from different modality combinations.

2017). In contrast to unimodal approaches, MSA has therefore become a cornerstone task in affective computing (Yu et al., 2024; Yang et al., 2025a; Zhao et al., 2025a; Ahuja et al., 2025).

Despite strong progress, most MSA models still rely on a fixed modality interaction pattern. A common line of work either treats all modalities as equally informative or assumes a dominant anchor modality (often text) to guide fusion. However, in real-world scenarios, an important challenge is that modality interaction pattern often shows great dynamics. Specifically, the most reliable modality (or modality combination) can vary substantially

056 across samples, and cross-modal conflicts can be
057 hierarchical and complicated. For example, the
058 combination of a smiling expression and positive
059 words delivers the positive signal, whereas audio
060 represents sarcasm (Bao et al., 2025). In this case,
061 multimodal interaction patterns with equal treat-
062 ment or over-reliance on unimodal are both failing
063 to resolve cross-modal polarity conflicts.

064 To further explore the existence and significance
065 of this challenge, we conducted a pilot study. As
066 shown in Fig. 1, three key findings emerge: (1)
067 *Trimodal fusion is not universally optimal*: predic-
068 tions from full-modality fusion can deviate substan-
069 tially from ground truth, even exhibiting opposite
070 polarity (e.g., -1.59 vs. 0.8). (2) *No single modal-
071 ity consistently dominates*: in the upper sample,
072 the audio-visual combination (0.31) outperforms
073 all text-involved combinations, contradicting the
074 common assumption of text dominance. (3) *Opti-
075 mal combinations are sample-dependent*: the up-
076 per sample favors bimodal fusion while the lower
077 benefits from unimodal input. Moreover, within
078 a single sample, different combinations can yield
079 conflicting polarities (-0.54 vs. 0.46), making in-
080 discriminate fusion unreliable.

081 These observations motivate a central question:
082 **how can we dynamically identify and leverage**
083 **sample-level preferences over modality com-**
084 **binations?** We decompose this question into
085 two subproblems. ① **Combination matching**:
086 for a given sample, which modalities should be
087 fused together and which should be kept separate?
088 ② **Fusion validation**: among the activated uni-
089 modal/bimodal/trimodal fusions, which fused rep-
090 resentations are reliable for final prediction?

091 To address the above question, we pro-
092 pose **Customized-Allocation Mixture-of-Experts**
093 **(CA-MoE)**, a novel framework that enables dy-
094 namic, sample-adaptive modality routing within
095 the MoE architecture. Specifically, CA-MoE con-
096 sists of two core complementary components:①
097 **Affinity-guided Customized Modality Alloca-**
098 **tion (ACMA)** acts as a *distributor* and solves
099 combination matching by routing modality com-
100 binations through an attraction–repulsion mecha-
101 nism, guided by Geometry–Gradient Affinity (G^2 -
102 Affinity), which jointly reflects semantic alignment
103 and optimization compatibility. ② **Reliability-**
104 **Aware Expert Selection (RAES)** is proposed as
105 a *validator*, which comprehensively evaluates the
106 reliability from the sentiment angular-prototype
107 proximity and competitive magnitude intensity of

108 representation, producing a reliability weighting
109 over experts for the final decision.

110 In summary, CA-MoE explicitly and dynami-
111 cally models sample-level dependency preferences
112 over modality combinations and prioritizes reliable
113 fused representations to improve sentiment predic-
114 tion. Our main contributions are three-folds:

- We identify and analyze an underexplored yet
115 challenging problem in MSA: How to dynami-
116 cally identify and leverage sample-level depen-
117 dency preference among different modality com-
118 binations? 119
- We propose CA-MoE featuring ACMA for iden-
120 tifying and routing sample-varied modal combi-
121 nations to expert and RAES for validating the
122 reliability after-fused expert representation to en-
123 courage most reliable experts contribute to the
124 final sentiment decision. 125
- Extensive experiments on three widely-used
126 MSA datasets demonstrate state-of-the-art per-
127 formance. Comprehensive ablations validate the
128 contribution of each proposed component. 129

130 2 Related Work

131 Most existing MSA approaches based on assump-
132 tions on modality priors can be categorized into two
133 groups: equally informative assignment methods
134 and dominant modality assignment methods. 134

135 **Equally informative assignment methods.** Pre-
136 vious works in this group treat all modalities’ in-
137 formation as a priori and pursue performance gains
138 through capturing richer semantic understanding
139 by designing high-capacity encoders (He et al.,
140 2025; Zhuang et al., 2024) or optimizing repre-
141 sentation strategy (Yu et al., 2021a; Hazarika et al.,
142 2020). Moreover, many studies also focus on cross-
143 modal interaction mechanisms. Early approaches
144 (Zadeh et al., 2017; Liu et al., 2018) typically advo-
145 cate for explicit interaction mechanisms for cross-
146 modal dynamic dependencies. Compared with their
147 static fusion paradigm, recent advances have turned
148 to cross-modal Transformer architectures, which
149 leverage cross-attention mechanisms for sample-
150 level interactions in a dynamic manner (Huang
151 et al., 2024; Xu et al., 2025; Gong et al., 2026;
152 Zhao et al., 2025b).

153 **Dominant modality assignment methods.**
154 These methods emphasize semantic guidance from
155 a dominant modality, usually text (Yang et al.,

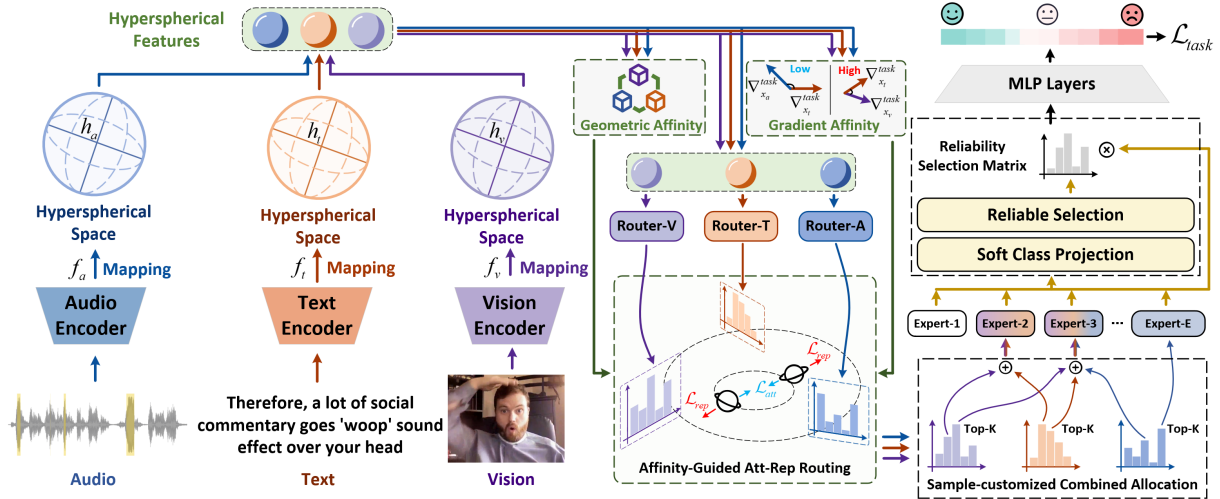


Figure 2: The overall architecture of CA-MoE. The blue, orange, and purple represent the audio, text, and vision modalities. Especially, we utilize colorless to indicate that there is no modal index to this expert. Simultaneously, various combinations of colors and quantities represent different modal combinations. For example, Expert-2 with two different color means the bimodal combination with text and vision.

2025b). They treat the dominant modality as an anchor and align auxiliary modalities through text-centric interaction mechanisms (Wang et al., 2025; Wu et al., 2025a; Yang et al., 2025c). However, recent studies have noted the issue of semantic noise introduced by text modality as a fixed prior. Some researches enable the model to dynamically select the dominant modality and adjust the contribution of each modality through additional unimodal modal supervision (Feng et al., 2024), sample gradient (Li et al., 2024), and routing mechanism (Gao et al., 2024; Fang et al., 2025).

Despite the promising results achieved for MSA, the above approaches ignore the possibility of modality combination as dominance, and these methods share an underlying limitation: their modality interaction patterns are usually fixed multimodal a priori and applied uniformly across all samples. However, in real-world applications, such an interaction pattern is often highly dynamic and even contains complicated, hierarchical intermodal conflicts. To solve this, we propose a dynamic modalities allocation approach to flexibly customize modal combinations and interaction. Furthermore, we prioritize reliability after-fused representation, thus improving predictions.

3 Methodology

3.1 Overall Architecture

As shown in Fig 2, which shows the overall workflow of CA-MoE. Specifically, we first ex-

tracts unimodal representations with corresponding encoders. The unimodal representations are mapped into hyperspherical space to obtain hyperspherical features for emotional polarity expression (Sec3.3). Then, Affinity-guided Customized Modality Allocation (ACMA) leverages Geometry-Gradient Affinity as a guide to allocate different modal combinations (e.g., unimodal, bimodal, or trimodal) to expert by attraction and repulsion routing (Sec3.4). Next, for the different after-fused results, Reliability-Aware Expert Selection (RAES) establishes emotion prototype proximity estimator to conduct comprehensive evaluation of angular-confidence and competitive intensity for reliability selection. Finally, the multimodal representation with the weights of reliability selection matrix is performed to boost MSA performance (Sec3.5).

3.2 Problem definition

Following with prior research (Zeng et al., 2021; Yu et al., 2021b) in MSA task, we aim to predict the sentiments intensity in videos by utilizing multimodal inputs. For the given utterance, the inputs data consists of text (t), vision (v) and audio (a) modalities. The sequences of three modalities are represented as, $X_v \in \mathbb{R}^{T_v \times d_v}$, $X_t \in \mathbb{R}^{T_t \times d_t}$, and $X_a \in \mathbb{R}^{T_a \times d_a}$, where $T_m, m \in t, v, a$ is the sequence length, and d_m corresponds to the embedding dimension of each modality. The objective is to predict continuous sentiment intensity values $y \in \mathbb{R}$ or a predefined set of C categories $y \in \mathbb{R}^C$.

3.3 Feature extraction and Mapping

We effectively encode the inputs with each modality $X_{m \in \{t, v, a\}}$ to obtain semantic information $F_m \in \mathbb{R}^{T_m \times d_m}$. Specifically, the BERT (Devlin et al., 2019) is used to extract text modality features F_t and the stacked transformer encoder layers [23] are leveraged to obtain the vision and audio modalities representations $F_{m \in \{v, a\}}$:

$$\begin{aligned} F_t &= \Phi_{BERT}(X_t) \in \mathbb{R}^{d_t} \\ F_m &= \Phi_{Trans}(X_m) \in \mathbb{R}^{d_m} \end{aligned} \quad (1)$$

Compared to Euclidean space, Hyperspherical embedding mitigates modality imbalance arising from feature magnitude discrepancies in similarity computation, owing to its boundedness and directional sensitivity (Ennajari et al., 2022). Normalizing features onto the unit hypersphere promotes cross-modal semantic alignment for emotional polarity modeling (Rizkallah et al., 2020; Li et al., 2025). The hypersphere space provides a scale invariant common semantic manifold for three modalities F_m , enabling framework to more reliably capture emotional polarity features. The process of mapping into hyperspherical features $h_m \in \mathbb{R}^{B \times d_s}$ by L2-normalizing is as follows:

$$h_m = \frac{F_m}{\|F_m\|_2 + \epsilon}, m \in \{t, v, a\} \quad (2)$$

where $\epsilon = 10^{-8}$ is a small constant for numerical stability. B and d_s are denoted as batch-size and hyperspherical space dimension, respectively

3.4 Affinity-guided Customized Modality Allocation (ACMA)

Geometry-Gradient Affinity (G²-Affinity) To resolve the first problem of combination matching, a comprehensive measurement basis should be designed to guide which modalities should be attracted together and which modalities should be repelled. Thus, we propose the G²-Affinity to commands dynamic routing of modal allocation. A high-quality modality combination should exhibit the two complementary aspects of strong gradient alignment and emotional semantic alignment. Firstly, to capture cross-modal affective consensus, we define a geometric affinity $A_{Ge}^{m,n}$ in the hyperspherical embedding space, where the cosine similarity between modalities directly reflects their emotional semantics directional coherence:

$$A_{Ge}^{m,n} = \frac{1}{2} \left(1 + h_m h_n^\top \right), \forall m, n \in \{t, v, a\}, m \neq n \quad (3)$$

Although geometric affinity $A_{Ge}^{m,n}$ effectively promotes alignment on emotional semantic level, it cannot reveal whether modalities are collaboratively optimized during training. Inspired by monitoring on gradient (Yu et al., 2020a; Wu et al., 2025b; Borsani et al., 2025), to bridge this gap, we further propose gradient affinity $A_{Gr}^{m,n}$ to capture gradient alignment for alleviating the potential gradient conflicts that may undermine learning stability and fusion efficacy:

$$\begin{aligned} \nabla_{h_m}^{\text{task}} &= \frac{\partial \mathcal{L}_{\text{task}}}{\partial h_m} \\ A_{Gr}^{m,n} &= \frac{1}{2} \left(1 + \cos \left(\nabla_{h_m}^{\text{task}}, \nabla_{h_n}^{\text{task}} \right) \right), \quad m \neq n \end{aligned} \quad (4)$$

G²-Affinity can be calculate by weighted average above two complementary component:

$$A_{m,n} = \beta A_{Ge}^{m,n} + (1 - \beta) A_{Gr}^{m,n} \quad (5)$$

where $\beta \in [0, 1]$ is a hyperparameter that controls the weights of geometric affinity and gradient affinity. Higher $A_{m,n}$ represents modality combinations that maintains semantic-level alignment and optimization-level harmony.

Affinity-Guided Attraction and Repulsion

Based on the guidance of G²-Affinity, the routing mechanism further acts on distributor to allocate expert with sample-customized combination. Formally, the CA-MoE consists of multiple experts, denoted as $e_1, e_2, \dots, e_{|E|}$, where $|E|$ represents the total number of experts. In order to obtain sample-customized modality combination, we leverage an independent router $R_{e,m}, e \in \{1, 2, \dots, |E|\}$ with corresponding learnable linear matrix $w_{e,m}$ for each modal m . Specifically, for the given embeddings h_m , the router engages the top-k-scores experts after $w_{e,m}$ and softmax operator. This process can be described as follows:

$$\begin{aligned} q_{e,m} &= \text{softmax}(w_{e,m}(h_m)), \quad \sum_{|E|} q_{e,m} = 1 \\ R_{e,m} &= \text{Top-k}(q_{e,m}, k) \end{aligned} \quad (6)$$

For the challenge of dynamically adapting to sample-level dependency preference, the routing mechanism is required to equip with the ability to customize the allocation of modal combinations. Specifically, high-affinity modality pairs share overlapping experts and low-affinity pairs activate disjoint expert subsets. Hence, we regard the router as a situational distributor, which appropriately modulates modalities to attract or repel.

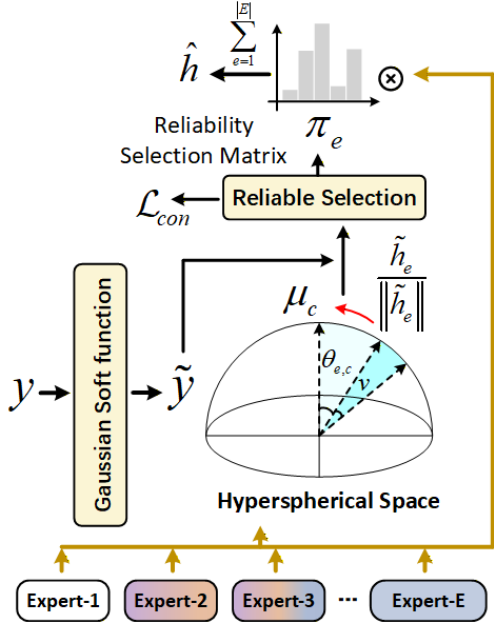


Figure 3: The details of the proposed RAES.

To achieve the allocation of similar expert subsets under attracted situation, the attraction Routing Loss \mathcal{L}_{att} is designed as controller to minimize the distance among high-affinity modal inputs q_m :

$$\mathcal{L}_{att} = \sum_{m,n} A_{m,n} [D_{KL}(q_m || q_n) + D_{KL}(q_n || q_m)] \quad (7)$$

Conversely, for the repelled situation, the Repulsion Routing Loss \mathcal{L}_{rep} controls router to disperse the expert selection involving potentially conflicting modalities by maximizing divergence among low-affinity modal inputs q_m :

$$\mathcal{L}_{rep} = - \sum_{m,n} (1 - A_{m,n}) D_{JS}(q_m || q_n) \quad (8)$$

where $m, n \in \{t, v, a\}, m \neq n$. The $D_{KL}(\cdot || \cdot)$ and $D_{JS}(\cdot || \cdot)$ are Kullback-Leibler and Jensen-Shannon divergence, respectively.

Through situational distributor controlled by \mathcal{L}_{att} and \mathcal{L}_{rep} , we achieve adaptive dependency preference allocation on sample-level by the guidance of affinity as $\tilde{h}_e \in \mathbb{R}^{B \times d_s}$:

$$\tilde{h}_e = \sum_m R_{e,m} h_m \quad (9)$$

3.5 Reliability-Aware Expert Selection (RAES)

Although different modal combinations have been assigned to various experts, determining the reli-

bility of after-fused modality features and converting them into contributions to the final task remains the hanging challenge. To this end, we introduce the RAES, as shown in Fig 3, that evaluates and refine expert outputs before final prediction.

The hyperspherical space provides a geometrically appropriate embedding space for MSA task, which is conducive to angular distance polarity discrimination and enables the marginal loss (Deng et al., 2019, 2020) to sharpen class boundaries and improve discriminative robustness.

From the angular-based perspective, our establishment of emotion prototype proximity estimator begins with a learnable category center μ_c with the label $y \in \mathbb{R}^C$. In hypersphere space, expert proximity scores $p_e(c)$ are formulated as:

$$\cos \theta_{e,c} = \mu_c^\top \frac{\tilde{h}_e}{\|\tilde{h}_e\|_2} \quad (10)$$

$$z_{e,c} = \begin{cases} s \cdot \cos(\theta_{e,c} + \nu), & c = y \\ s \cdot \cos(\theta_{e,c}), & c \neq y \end{cases}$$

$$p_e(c) = \text{softmax}_c(z_{e,\cdot})$$

where s is the hyperparameter scaled to a numerical scale suitable for optimization. And ν is the fixed angle margin.

To adapt to the calculation of $p_e(c)$, we project discrete integer labels into a Gaussian soft target distribution with $\tilde{y} = \text{softmax}\left(\exp\left(-\frac{(c-y)^2}{2\sigma^2}\right)\right)$. Then, the expert angular-confidence calculated as:

$$\varpi_e = \sum_c p_e(c) \cdot \tilde{y} \quad (11)$$

The form of expert proximity scores $K_e = -\sum_c p_e(c) \log p_e(c)$ could measure the directional consensus of after-fused outputs from different experts. $H_e = \|\tilde{h}_e\|_2$ could further complementarily reflect the activation and competitive intensity among experts. Thus, reliability includes high angular agreement and magnitude responsiveness. And the final reliability selection matrix and weighted representations are defined as:

$$\pi_e = \alpha \cdot \text{softmax}(\gamma_1 H_e - \gamma_2 K_e) + (1 - \alpha) \varpi_e$$

$$\hat{h} = \sum_{e=1}^{|E|} \text{softmax}(\pi_e) \tilde{h}_e \quad (12)$$

where γ_1 and γ_2 are hyperparameters to control the weight on the components of expert ability. The α weights the importance of expert confidence and ability in reliability selection matrix.

Table 1: Performance comparison between the CA-MoE and other baselines on the CMU-MOSI and CMU-MOSEI. The previous SOTA values are annotated with underscores, while the current SOTA values are highlighted in bold. In Acc-2 and F1-Score, the values on the left side of “/” represent “negative/non-negative”, while the values on the right represent “negative/positive”. A lower MAE value indicates better performance, while higher values of the other evaluation metrics (Corr, ACC-2, F1-score, ACC-7) signify superior performance.

Model	Date	CMU-MOSI					CMU-MOSEI				
		MAE↓	Corr↑	ACC-2↑	F1-score↑	ACC-7↑	MAE↓	Corr↑	ACC-2↑	F1-score↑	ACC-7↑
TFN	EMNLP,2017	0.901	0.698	-80.8	-80.7	34.9	0.593	0.700	-82.5	-82.1	50.2
MuT	ACL,2019	0.861	0.711	81.5/84.1	80.6/83.9	-	0.580	0.703	-82.5	-82.3	-
ICCN	AAAI,2020	0.862	0.714	-83.1	-83.0	39.0	0.565	0.713	-84.2	-84.2	51.6
MISA	ACM,2020	0.804	0.764	80.8/82.1	80.8/82.0	42.3	0.568	0.724	82.6/84.2	82.7/84.0	-
Self-MM	AAAI,2021	0.712	0.795	82.5/84.8	82.7/84.9	45.8	0.529	0.767	82.7/85.0	83.0/84.9	53.5
TETFN	PR,2023	0.717	0.800	84.1/86.1	83.8/86.1	-	0.551	0.748	84.3/85.2	84.2/85.3	-
SKEAFN	IF,2023	0.740	0.784	84.4/86.4	84.5/86.7	45.2	0.540	0.763	83.7/86.0	83.4/86.1	52.8
CRNet	KBS,2024	0.712	0.797	-86.4	-86.4	47.4	0.541	0.771	-86.2	-86.1	53.8
FGTI	ICASSP,2024	0.702	0.791	-85.8	-85.8	48.0	0.536	0.771	-86.0	-86.0	53.4
DTN	IF,2024	0.714	0.807	-86.2	-86.2	48.1	0.579	0.788	-86.3	-86.3	52.5
EMT	TAC,2023	0.705	0.798	83.3/85.0	83.2/85.0	47.4	0.527	0.774	83.4/86.0	83.7/86.0	54.5
TIEMFF	PR,2024	0.727	0.698	-85.7	-85.8	-	0.548	0.705	-85.9	-81.0	-
DLF	AAAI,2025	0.731	0.781	-85.1	-85.0	47.1	0.536	0.764	-85.4	-85.3	53.9
EMOE	CVPR,2025	0.725	0.792	83.5/85.2	83.4/85.2	45.8	0.543	0.759	82.2/85.1	82.4/84.9	53.1
MMLN	ESWA,2026	0.707	0.797	83.3/86.3	83.6/86.2	47.8	0.530	0.773	85.3/86.7	85.3/86.8	54.6
CR-GAC	ESWA,2026	0.684	0.814	85.6/87.5	85.6/87.5	44.3	0.521	0.793	84.4/86.4	84.4/86.5	54.7
CA-MoE(ours)		0.679	0.822	86.0/88.4	86.2/88.5	48.5	0.515	0.802	85.8/87.5	83.7/87.3	55.4
Δ SOTA		↓0.005	↑0.008	↑0.4/↑0.9	↑0.6/↑1.0	↑0.4	↓0.006	↑0.009	↑0.5/↑0.7	↓1.6/↑0.5	↑0.7

3.6 Output and learning Objectives

We put the weighted representations into MLP layers for linear projection:

$$\hat{y} = \text{MLP}(\hat{h}) \quad (13)$$

The model is optimized by the total loss:

$$\mathcal{L}_{total} = \mathcal{L}_{task} + \lambda_1 \mathcal{L}_{att} + \lambda_2 \mathcal{L}_{rep} + \lambda_3 \mathcal{L}_{con} + \lambda_4 \mathcal{L}_{bal} \quad (14)$$

where \mathcal{L}_{task} represents primary task loss (Hazari et al., 2020). Expert confidence loss $\mathcal{L}_{con} = D_{KL}(p_e(c) || \hat{y})$ is used to optimize the category center. \mathcal{L}_{bal} is introduced to mitigate expert underutilization (Fedus et al., 2022). And $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ control contribution of each loss term.

4 Experiments

4.1 Experimental Settings

Datasets. We evaluate CA-MoE on three benchmark MSA datasets: CMU-MOSI (Zadeh et al., 2016), CMU-MOSEI (Zadeh et al., 2018), and CH-SIMS (Yu et al., 2020b). The detailed statistics are summarized in Appendix A.

Implementation Details. Following standard protocols, we adopt MAE, Corr, Acc-2, F1, and Acc-7 as evaluation metrics. All experiments are

conducted on an NVIDIA RTX 3060 GPU. The Appendix B shows detailed implementation settings. We also provide efficiency analysis (Appendix C) and algorithm process (Appendix D).

4.2 Comparison with State-of-the-Art

We compare CA-MoE with competitive baselines, including TFN (Zadeh et al., 2017), MuT (Tsai et al., 2019), MISA (Hazari et al., 2020), and recent SOTA methods like MMLN (Li et al., 2026) and CR-GAC (Chen et al., 2026). The quantitative results on CMU-MOSI and CMU-MOSEI are reported in Table 1, and results on CH-SIMS are shown in Table 2.

Results on CMU-MOSI and CMU-MOSEI.

As observed in Table 1, CA-MoE consistently outperforms existing methods across most metrics. Notably, on the CMU-MOSI dataset, our method achieves a new state-of-the-art F1-score of 88.5% and Acc-2 of 88.4%, surpassing CR-GAC by 1.0% and 0.9%, respectively. Furthermore, CA-MoE achieves the lowest MAE (0.679) and highest correlation (0.822), demonstrating its capability to learn precise sentiment intensity. Similar improvements are observed on CMU-MOSEI, where CA-MoE maintains competitive performance, validating its robustness on large-scale data.

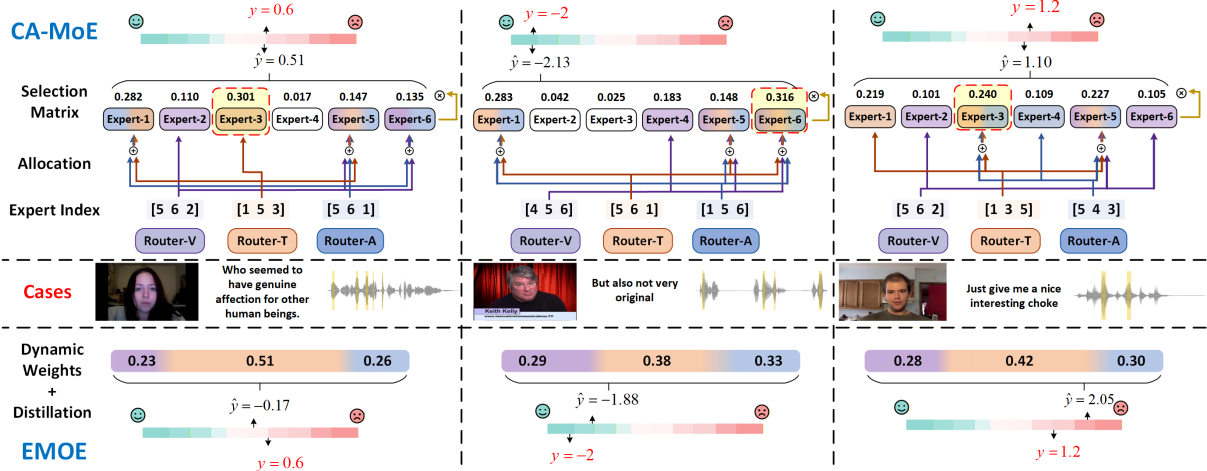


Figure 4: Qualitative comparison between CA-MoE and EMOE. The top panel visualizes the ACMA and RAES of CA-MoE, where red dashed boxes highlight high-reliability experts. The bottom panel shows EMOE’s methods. CA-MoE achieves more accurate predictions by customizing expert combinations and selecting reliable presentations.

Table 2: Performance comparison on **CH-SIMS** datasets. The optimal results are highlighted in bold.

Model	MAE	Corr	Acc-2	Acc-3	Acc-5	F1
TFN	0.432	0.591	78.3	65.1	39.3	78.6
LMF	0.441	0.575	77.7	64.6	40.5	77.8
MuT	0.442	0.581	78.2	65.7	40.0	78.5
Self-MM	0.411	0.601	78.6	66.1	43.1	78.6
CENet	0.470	0.539	77.9	62.5	33.9	77.5
TETFN	0.420	0.577	81.1	63.2	41.7	80.2
EMOE	0.449	0.519	75.7	62.1	40.0	75.9
DLF	0.446	0.563	80.1	64.3	39.8	78.8
MMLN	<u>0.406</u>	0.602	<u>80.7</u>	<u>67.1</u>	<u>44.6</u>	<u>80.9</u>
CR-GAC	0.412	<u>0.613</u>	<u>80.7</u>	66.3	43.9	<u>80.9</u>
CA-MoE(ours)	0.398	0.624	81.5	68.0	45.5	81.6
Δ SOTA	\downarrow 0.008	\uparrow 0.011	\uparrow 0.8	\uparrow 0.9	\uparrow 0.9	\uparrow 0.7

Results on CH-SIMS. To verify the generalization ability on different languages, we report results on the Chinese dataset CH-SIMS in Table 2. CA-MoE achieves the best performance in MAE (0.398) and Corr (0.624), significantly outperforming the runner-up MMLN (MAE 0.406). Moreover, our method demonstrates substantial gains in multi-class classification, improving Acc-5 and Acc-3 by 0.9% compared to the previous best results.

4.3 Qualitative Analysis

To intuitively demonstrate the advantages of CA-MoE in **sample-level dependency preferences**, Figure 4 compares its performance with the baseline model EMOE on three representative samples. Although methods such as EMOE achieve flexible modal interaction through dynamic weight mechanisms and can effectively utilize all available modal

information, their fusion strategies essentially still adopt a **fixed interaction paradigm**, making it difficult to precisely adapt to the modal dependency characteristics of different samples, particularly as in Cases 1 and 3. In contrast, CA-MoE allocates modality combinations under affinity guidance and further prioritizes high-scoring representations. CA-MoE identifies and prioritizes key experts (e.g., Expert-3 highlighted in the red dashed box) while avoiding falling into the fixed paradigm, thereby rectifying polarity errors and outputting predictions closer to the ground truth.

4.4 Ablation Study

Impact of Key Components. To verify the effectiveness of the model components, we conducted ablation studies on the CMU-MOSI dataset (in Table 3). The results indicate that removing any module leads to a decline in performance. Most notably, the impact of RAES as the **validator** is significant; its removal causes the MAE to rise from 0.679 to 0.760 and Acc-7 to drop by 4.8%, demonstrating the core role of reliability assessment. Furthermore, the exclusion of \mathcal{L}_{att} or \mathcal{L}_{rep} also results in decreased metrics, confirming the necessity of G^2 -Affinity in guiding modality routing.

Visualization of feature distribution. To validate the impact of the Reliability-Aware Expert Selection (RAES) module, we visualize the feature distributions of the ablation variant (w/o RAES) and the full CA-MoE in (in Fig. 5). As observed, the full CA-MoE exhibits a distinct, continuous manifold structure where samples transition

Table 3: **Ablation study of different components on CMU-MOSI dataset.** “w/o” denotes the removal of a specific module from the full model. The best results are highlighted in bold.

Model	MAE↓	Corr↑	Acc-2(%)↑	F1-score(%)↑	Acc-7(%)↑
CA-MoE	0.679	0.822	86.0/88.4	86.2/88.5	48.8
w/o \mathcal{L}_{att}	0.720	0.800	84.2/86.7	84.4/86.8	46.0
w/o \mathcal{L}_{rep}	0.738	0.789	83.2/84.9	83.0/85.1	43.8
w/o \mathcal{L}_{con}	0.713	0.803	84.8/85.7	85.0/87.7	47.1
w/o RAES	0.760	0.772	82.0/83.7	82.1/83.7	44.0

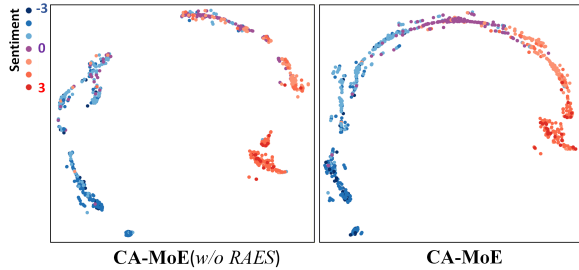
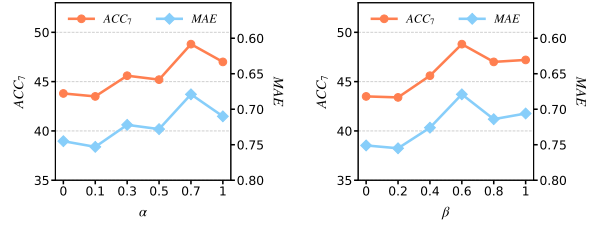


Figure 5: t-SNE visualization of feature embeddings on the CMU-MOSI test set.

smoothly from negative (blue) to positive (red), aligning perfectly with the hyperspherical feature space design. In contrast, the distribution without RAES appears fragmented and disordered. This demonstrates that RAES effectively enhances the model’s ability to learn discriminative and continuous sentiment representations.

Parameter Sensitivity Analysis. We further demonstrate that the model exhibits good robustness to hyper-parameters α and β (in Fig. 6). Experiments show that performance remains stable for α within the range of $[0.3, 0.7]$, peaking at 0.7, while β performs best around 0.6. This suggests that the model does not rely on specific parameter points but displays strong generalization stability when balancing auxiliary loss weights and the ratio of geometric to gradient affinity.

Impact of Expert Configuration. Regarding the expert settings (in Table 4), the investigation reveals that a total expert count of $|E| = 6$ serves as the optimal balance point for model capacity and generalization; too few experts lead to insufficient representation, while too many introduce redundancy. For the selection strategy, $k = 3$ proves to be the optimal solution. Compared to a single path ($k = 1$) or full activation ($k = 6$), this setting effectively leverages the complementarity of modality combinations while maintaining the advantages of sparse computation.



(a) Parameter Analysis on α (b) Parameter Analysis on β

Figure 6: **Parameter Sensitivity Analysis.** We investigate the impact of hyper-parameters α and β on the CMU-MOSI dataset. The results indicate that the model is robust within a reasonable range.

Table 4: **Ablation Study on CMU-MOSI Dataset.** We investigate the impact of the number of experts ($|E|$) and the Top- k selection strategy. The results are reported in terms of MAE, Correlation (Corr), Acc-2, F1-score, and Acc-7. The default setting is $|E| = 6$ and $k = 3$.

Setting	MAE↓	Corr↑	Acc-2(%)↑	F1-score(%)↑	Acc-7(%)↑
<i>Impact of Expert Number ($k = 3$)</i>					
$ E = 4$	0.716	0.806	84.3/85.2	84.1/85.5	45.7
$ E = 6$	0.679	0.822	86.0/88.4	86.2/88.5	48.8
$ E = 8$	0.704	0.810	84.5/86.1	85.2/86.6	46.2
$ E = 10$	0.742	0.791	83.0/85.2	82.9/84.5	44.1
$ E = 12$	0.767	0.782	82.7/84.6	82.7/84.4	43.8
$ E = 16$	0.762	0.776	82.2/84.0	82.5/84.1	43.7
<i>Impact of Top-k Selection ($E = 6$)</i>					
$k = 1$	0.721	0.811	83.8/85.7	83.7/85.2	46.2
$k = 2$	0.727	0.802	83.2/85.3	83.4/85.2	45.8
$k = 3$	0.679	0.822	86.0/88.4	86.2/88.5	48.8
$k = 4$	0.701	0.817	85.5/87.7	85.3/87.2	47.5
$k = 5$	0.692	0.819	85.8/87.9	86.0/88.2	48.0
$k = 6$	0.705	0.808	85.2/87.3	85.2/87.7	47.7

5 Conclusion

This paper explores the core issue of *How to dynamically identify and leverage sample-level dependency preferences among different modality combinations?* in the MSA task. Furthermore, we decompose it into two subproblems: combination matching and fusion validation. To this end, we propose a sample-level and dynamical MoE-based method, the Customized-Allocation Mixture-of-Experts (CA-MoE), which is designed for reliable modalities combination allocation. First, CA-MoE utilizes Geometry-Gradient Affinity as the attraction–repulsion routing guidance to solve customized-combination matching. Then, the validator evaluates the reliability of the fused representation for yielding a weighted final prediction. The effectiveness of CA-MoE has been thoroughly validated against MAS task datasets. We hope this work will provide inspiration for the future research on multimodal interaction paradigm.

514 Limitations

515 Although CA-MoE achieves significant perfor-
516 mance improvements and computational efficiency,
517 there are still some limitations that need to be ad-
518 dressed in future work. Firstly, due to constraints in
519 computational resources, we did not re-implement
520 the baseline methods but directly referenced the
521 experimental results reported in their original pa-
522 pers. To ensure the fairness of the comparison, we
523 strictly adopted the exact same dataset splits, fea-
524 ture extraction protocols, and evaluation metrics as
525 those of prior studies. Secondly, while this paper
526 effectively resolves cross-modal conflicts via the
527 affinity-guided routing mechanism, the model cur-
528 rently relies on pre-extracted features (e.g., BERT).
529 In real-world environments with high background
530 noise or occlusion, the quality of these unimodal
531 features may degrade, potentially affecting the cal-
532 culation of geometric and gradient affinities. This
533 imposes higher demands on feature robustness and
534 offers opportunities for future optimization. Finally,
535 the proposed routing mechanism implicitly learns
536 sample-specific modality preferences through con-
537 strained modality distributions. While effective,
538 this paradigm remains latent; future work should
539 aim to explicitly model and interpret the underlying
540 criteria governing modality combination decisions.

541 References

542 Garvit Ahuja, Alireza Alaei, and Umapada Pal. 2025. A
543 new multimodal sentiment analysis for images con-
544 taining textual information. *Multimedia tools and*
545 *applications*, 84(21):23745–23774.

546 Tadas Baltruaitis, Chaitanya Ahuja, and Louis-Philippe
547 Morency. 2019. Multimodal machine learning: A
548 survey and taxonomy. *IEEE transactions on pattern*
549 *analysis and machine intelligence*, 41(2):423–443.

550 Yongtang Bao, Xin Zhao, Peng Zhang, Yue Qi, and
551 Haojie Li. 2025. Hian: A hybrid interactive attention
552 network for multimodal sarcasm detection. *Pattern*
553 *Recognition*, 164:111535.

554 Thomas Borsani, Andrea Rosani, Giuseppe Nicosia, and
555 Giuseppe Di Fatta. 2025. Gradient similarity surgery
556 in multi-task deep learning. In *Joint European Con-*
557 *ference on Machine Learning and Knowledge Dis-*
558 *covery in Databases*, pages 95–111. Springer.

559 Haoran Chen, Jiapeng Liu, Zuhe Li, Yushan Pan, Hong-
560 wei Tao, Huaiguang Wu, Yunyang Wang, and Chen-
561 guang Yang. 2026. Cr-gac: Cross-modal recombi-
562 nation via graph-attention collaborative optimization
563 for multimodal sentiment analysis. *Expert Systems*
564 *with Applications*, 298:129805.

Jiankang Deng, Jia Guo, Tongliang Liu, Mingming
Gong, and Stefanos Zafeiriou. 2020. Sub-center arc-
face: Boosting face recognition by large-scale noisy
web faces. In *European Conference on Computer*
Vision, pages 741–757. Springer. 565 566 567 568 569

Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos
Zafeiriou. 2019. Arcface: Additive angular margin
loss for deep face recognition. In *Proceedings of*
the IEEE/CVF conference on computer vision and
pattern recognition, pages 4690–4699. 570 571 572 573 574

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2019. Bert: Pre-training of deep
bidirectional transformers for language understand-
ing. In *Proceedings of the 2019 conference of the*
North American chapter of the association for com-
putational linguistics: human language technologies,
volume 1 (long and short papers), pages 4171–4186. 575 576 577 578 579 580 581

Hafsa Ennajari, Nizar Bouguila, and Jamal Bentahar.
2022. Knowledge-enhanced spherical representation
learning for text classification. In *Proceedings of the*
2022 SIAM international conference on data mining
(SDM), pages 639–647. SIAM. 582 583 584 585 586

Yiyang Fang, Wenke Huang, Guancheng Wan, Kehua
Su, and Mang Ye. 2025. Emoe: Modality-specific
enhanced dynamic emotion experts. In *Proceedings*
of the Computer Vision and Pattern Recognition Con-
ference, pages 14314–14324. 587 588 589 590 591

William Fedus, Barret Zoph, and Noam Shazeer. 2022.
Switch transformers: Scaling to trillion parameter
models with simple and efficient sparsity. *Journal of*
Machine Learning Research, 23(120):1–39. 592 593 594 595

Xinyu Feng, Yuming Lin, Lihua He, You Li, Liang
Chang, and Ya Zhou. 2024. Knowledge-guided dy-
namic modality attention fusion framework for multi-
modal sentiment analysis. In *Findings of the Associ-*
ation for Computational Linguistics: EMNLP 2024,
pages 14755–14766. 596 597 598 599 600 601

Zixian Gao, Disen Hu, Xun Jiang, Huimin Lu, Heng Tao
Shen, and Xing Xu. 2024. Enhanced experts with
uncertainty-aware routing for multimodal sentiment
analysis. In *Proceedings of the 32nd ACM Interna-*
tional Conference on Multimedia, pages 9650–9659. 602 603 604 605 606

Peizhu Gong, Jin Liu, Xiliang Zhang, Xingye Li, Lai
Wei, and Huihua He. 2026. Towards robust senti-
ment analysis with multimodal interaction graph
and hybrid contrastive learning. *Pattern Recognition*,
169:111870. 607 608 609 610 611

Devamanyu Hazarika, Roger Zimmermann, and Sou-
janya Poria. 2020. Misa: Modality-invariant and-
specific representations for multimodal sentiment
analysis. In *Proceedings of the 28th ACM interna-*
tional conference on multimedia, pages 1122–1131. 612 613 614 615 616

Xilin He, Haijian Liang, Boyi Peng, Weicheng Xie,
Muhammad Haris Khan, Siyang Song, and Zitong
Yu. 2025. Msamba: Exploring multimodal sentiment
analysis with state space models. In *Proceedings*
617 618 619 620

621	<i>of the AAAI Conference on Artificial Intelligence</i> ,	<i>of the AAAI Conference on Artificial Intelligence</i> ,	677
622	volume 39, pages 1309–1317.	volume 39, pages 1601–1609.	678
623	Jiehui Huang, Jun Zhou, Zhenchao Tang, Jiaying Lin, and Calvin Yu-Chian Chen. 2024.	Yongliang Wu, Shiji Zhou, Mingzhuo Yang, Lianzhe Wang, Heng Chang, Wenbo Zhu, Xinting Hu, Xiao Zhou, and Xu Yang. 2025b. Unlearning concepts in diffusion model via concept domain correction and concept preserving gradient. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> ,	679
624	Tmbl: Transformer-based multimodal binding learning model for multimodal sentiment analysis.	volume 39, pages 8496–8504.	680
625	<i>Knowledge-Based Systems</i> , 285:111346.	Qinfu Xu, Yiwei Wei, Chunlei Wu, Leiquan Wang, Shaozu Yuan, Jie Wu, Jing Lu, and Hengyang Zhou. 2025. Towards multimodal sentiment analysis via hierarchical correlation modeling with semantic distribution constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> ,	681
626		volume 39, pages 21788–21796.	682
627		Hao Yang, Yanyan Zhao, Yang Wu, Shilong Wang, Tian Zheng, Hongbo Zhang, Zongyang Ma, Wanxiang Che, Shijin Wang, Si Wei, and 1 others. 2025a. Large language models meet text-centric multimodal sentiment analysis: A survey. <i>Science China Information Sciences</i> , 68(10):1–29.	683
628	Xiang Li, Zhiqiang Dong, Xianfu Cheng, Dezhuang Miao, Haijun Zhang, Tianbo Wang, Xiaoming Zhang, and Zhoujun Li. 2026. A multi-scale representation and multi-level decision learning network for multimodal sentiment analysis. <i>Expert Systems with Applications</i> , 297:129341.	Hao Yang, Yanyan Zhao, Yang Wu, Shilong Wang, Tian Zheng, Hongbo Zhang, Zongyang Ma, Wanxiang Che, Shijin Wang, Si Wei, and 1 others. 2025b. Large language models meet text-centric multimodal sentiment analysis: A survey. <i>Science China Information Sciences</i> , 68(10):1–29.	684
629		Yang Yang, Xunde Dong, and Yupeng Qiang. 2025c. Mse-adapter: A lightweight plugin endowing llms with the capability to perform multimodal sentiment analysis and emotion recognition. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> ,	685
630		volume 39, pages 25642–25650.	686
631		Hongqi Yu, Fei Tang, Lei Zhang, Randy Gomez, Eric Nichols, and Guangliang Li. 2024. Improving perceived emotional intelligence of embodied chatbot haru via multi-modal interaction. In <i>2024 IEEE International Conference on Robotics and Biomimetics (ROBIO)</i> , pages 51–58. IEEE.	687
632		Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. 2020a. Gradient surgery for multi-task learning. <i>Advances in neural information processing systems</i> , 33:5824–5836.	688
633		Wenmeng Yu, Hua Xu, Fanyang Meng, Yilin Zhu, Yixiao Ma, Jiele Wu, Jiyun Zou, and Kaicheng Yang. 2020b. Ch-sims: A chinese multimodal sentiment analysis dataset with fine-grained annotation of modality. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> ,	689
634	Xinyu Li, Wenqing Ye, Yueyi Zhang, and Xiaoyan Sun. 2024. Grace: Gradient-based active learning with curriculum enhancement for multimodal sentiment analysis. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 5702–5711.	Wenmeng Yu, Hua Xu, Ziqi Yuan, and Jiele Wu. 2021a. Learning modality-specific representations with self-supervised multi-task learning for multimodal sentiment analysis. In <i>Proceedings of the AAAI conference on artificial intelligence</i> ,	690
635		volume 35, pages 10790–10797.	691
636			692
637			693
638			694
639	Zuhe Li, Panbo Liu, Yushan Pan, Weiping Ding, Jun Yu, Haoran Chen, Weihua Liu, Yiming Luo, and Hao Wang. 2025. Multimodal sentiment analysis based on disentangled representation learning and cross-modal-context association mining. <i>Neurocomputing</i> , 617:128940.		695
640			696
641			697
642			698
643			699
644			700
645	Zhun Liu, Ying Shen, Varun Bharadhwaj Lakshminarasimhan, Paul Pu Liang, AmirAli Bagher Zadeh, and Louis-Philippe Morency. 2018. Efficient low-rank multimodal fusion with modality-specific factors. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2247–2256.		701
646			702
647			703
648			704
649			705
650			706
651			707
652	Soujanya Poria, Erik Cambria, Devamanyu Hazarika, Navonil Majumder, Amir Zadeh, and Louis-Philippe Morency. 2017. Context-dependent sentiment analysis in user-generated videos. In <i>Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: Long papers)</i> , pages 873–883.		708
653			709
654			710
655			711
656			712
657			713
658			714
659	Sandra Rizkallah, Amir F Atiya, and Samir Shaheen. 2020. A polarity capturing sphere for word to vector representation. <i>Applied Sciences</i> , 10(12):4386.		715
660			716
661			717
662	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>Proceedings of the conference. Association for computational linguistics. Meeting</i> , volume 2019, page 6558.		718
663			719
664			720
665			721
666			722
667			723
668			724
669			725
670			726
671			727
672			728
673	Pan Wang, Qiang Zhou, Yawen Wu, Tianlong Chen, and Jingtong Hu. 2025. Dlf: Disentangled-language-focused multimodal sentiment analysis. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> ,		729
674	volume 39, pages 21180–21188.		730
675			731
676			732
677			733
678			734
679			735
680			736
681			737
682			738
683			739
684			740
685			741
686			742
687			743
688			744
689			745
690			746
691			747
692			748
693			749
694			750
695			751
696			752
697			753
698			754
699			755
700			756
701			757
702			758
703			759
704			760
705			761
706			762
707			763
708			764
709			765
710			766
711			767
712			768
713			769
714			770
715			771
716			772
717			773
718			774
719			775
720			776
721			777
722			778
723			779
724			780
725			781
726			782
727			783
728			784
729			785
730			786
731			787
732			788
733			789
734			790
735			791
736			792
737			793
738			794
739			795
740			796
741			797
742			798
743			799
744			800
745			801
746			802
747			803
748			804
749			805
750			806
751			807
752			808
753			809
754			810
755			811
756			812
757			813
758			814
759			815
760			816
761			817
762			818
763			819
764			820
765			821
766			822
767			823
768			824
769			825
770			826
771			827
772			828
773			829
774			830
775			831
776			832
777			833
778			834
779			835
780			836
781			837
782			838
783			839
784			840
785			841
786			842
787			843
788			844
789			845
790			846
791			847
792			848
793			849
794			850
795			851
796			852
797			853
798			854
799			855
800			856
801			857
802			858
803			859
804			860
805			861
806			862
807			863
808			864
809			865
810			866
811			867
812			868
813			869
814			870
815			871
816			872
817			873
818			874
819			875
820			876
821			877
822			878
823			879
824			880
825			881
826			882
827			883
828			884
829			885
830			886
831			887
832			888
833			889
834			890
835			891
836			892
837			893
838			894
839			895
840			896
841			897
842			898
843			899
844			900
845			901
846			902
847			903
848			904
849			905
850			906
851			907
852			908
853			909
854			910
855			911
856			912
857			913
858			914
859			915
860			916
861			917
862			918
863			919
864			920
865			921
866			922
867			923
868			924
869			925
870			926
871			927
872			928
873			929
874			930
875			931
876			932
877			933
878			934
879			935
880			936
881			937
882			938
883			939
884			940
885			941
886			942
887			943
888			944
889			945
890			946
891			947
892			948
893			949
894			950
895			951
896			952
897			953
898			954
899			955
900			956
901			957
902			958
903			959
904			960
905			961
906			962
907			963
908			964
909			965
910			966
911			967
912			968
913			969
914			970
915			971
916			972
917			973
918			974
919			975
920			976
921			977
922			978
923			979
924			980
925			981
926			982
927			983
928			984
929			985

- 735 Wenmeng Yu, Hua Xu, Ziqi Yuan, and Jiele Wu. 2021b.
736 Learning modality-specific representations with self-
737 supervised multi-task learning for multimodal sen-
738 timent analysis. In *Proceedings of the AAAI con-
739 ference on artificial intelligence*, volume 35, pages
740 10790–10797.
- 741 Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cam-
742 bria, and Louis-Philippe Morency. 2017. Tensor
743 fusion network for multimodal sentiment analysis.
744 In *Proceedings of the 2017 Conference on Empiri-
745 cal Methods in Natural Language Processing*, pages
746 1103–1114.
- 747 Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-
748 Philippe Morency. 2016. Multimodal sentiment in-
749 tensity analysis in videos: Facial gestures and verbal
750 messages. *IEEE Intelligent Systems*, 31(6):82–88.
- 751 AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria,
752 Erik Cambria, and Louis-Philippe Morency. 2018.
753 Multimodal language analysis in the wild: Cmu-
754 mosei dataset and interpretable dynamic fusion graph.
755 In *Proceedings of the 56th Annual Meeting of the As-
756 sociation for Computational Linguistics (Volume 1:
757 Long Papers)*, pages 2236–2246.
- 758 Ying Zeng, Sijie Mai, and Haifeng Hu. 2021. Which
759 is making the contribution: Modulating unimodal
760 and cross-modal dynamics for multimodal sentiment
761 analysis. In *Findings of the Association for Computa-
762 tional Linguistics: EMNLP 2021*, pages 1262–1274.
- 763 Kai Zhao, Mingsheng Zheng, Qingguan Li, and Jianing
764 Liu. 2025a. Multimodal sentiment analysis-a com-
765 prehensive survey from a fusion methods perspective.
766 *IEEE Access*.
- 767 Xianbing Zhao, Xuejiao Li, Ronghuan Jiang, and
768 Buzhou Tang. 2025b. Decoupled cross-attribute cor-
769 relation network for multimodal sentiment analysis.
770 *Information Fusion*, 117:102897.
- 771 Yan Zhuang, Yanru Zhang, Zheng Hu, Xiaoyue Zhang,
772 Jiawen Deng, and Fuji Ren. 2024. Glomo: Global-
773 local modal fusion for multimodal sentiment analysis.
774 In *Proceedings of the 32nd ACM International Con-
775 ference on Multimedia*, pages 1800–1809.

A Datasets

We evaluate CA-MoE on three benchmark MSA datasets: CMU-MOSI, CMU-MOSEI, and CH-SIMS. The detailed statistics are summarized in Table 5. **CMU-MOSI** is a widely used dataset containing 2,199 video segments (1,284 for training, 299 for validation, and 686 for testing) annotated with sentiment intensity scores. **CMU-MOSEI** is a large-scale dataset consisting of 22,856 movie review clips, which are split into 16,326, 1,871, and 4,659 samples for training, validation, and testing, respectively. **CH-SIMS** is a Chinese multimodal sentiment analysis dataset containing 2,281 refined video segments with imbalanced modality annotations.

Table 5: Datasets details.

Dataset	Train	Valid	Test	All
CMU-MOSI	1284	299	686	2199
CMU-MOSEI	16326	1871	4659	22856
CH-SIMS	1368	456	457	2281

B Implementation Details.

Hyper-parameter	MOSI	MOSEI	SIMS
β	0.6	0.6	0.6
$ E $	6	8	6
$Top-k$	3	3	3
s, ν	30,0.3	30,0.3	30,0.3
$\alpha, \gamma_1, \gamma_2$	0.7,0.8,30	0.7,0.8,30	0.7,0.8,30
$\lambda_1, \lambda_2, \lambda_3, \lambda_3$	1,1,0.1,0.01	1,1,0.1,0.01	1,1,0.1,0.01
Batch size	64	16	32
Epoch	100	100	100
Optimizer	Adam	Adam	Adam
Learning rate	1e-4	1e-4	1e-4

Table 6: Hyper-parameters setting.

Following standard protocols, we adopt MAE, Pearson Correlation (Corr), Acc-2, F1-score, and Acc-7 as evaluation metrics. All experiments are implemented on the PyTorch framework using an NVIDIA RTX 3060 GPU. The models are trained for 100 epochs using the Adam optimizer with a learning rate of $1e^{-4}$. The batch sizes are set to 64, 16, and 32 for CMU-MOSI, CMU-MOSEI, and CH-SIMS, respectively. Key hyperparameters, such as the number of experts $|E|$ and Top- k selection, vary by dataset and are detailed in Table 6.

C Efficiency Analysis

Table 7 compares the computational overhead and performance of the proposed CA-MoE against the

state-of-the-art EMOE on the CMU-MOSI dataset. CA-MoE demonstrates superior efficiency, achieving a significantly faster runtime per epoch (16 s vs. 30 s) and requiring less than half the parameters (122 M vs. 317 M). Remarkably, this lightweight design facilitates better performance rather than compromising it; CA-MoE yields a lower Mean Absolute Error (0.679 vs. 0.725), proving that our model structure achieves an optimal balance between computational efficiency and predictive accuracy.

Table 7: Computational overhead.

Model	Parameters \downarrow	Time / Epoch \downarrow	MAE \downarrow
EMOE	317 M	30 s	0.725
CA-MoE	122 M	16 s	0.679

D Algorithm Process

In this work, we present the core computational pipeline of our proposed Customized-Allocation Mixture-of-Experts (CA-MoE) using pseudocode to enhance clarity and reproducibility. Specifically, Algorithm 1 delineates the full forward inference process—from raw multimodal inputs to customized modality allocation. Concurrently, Algorithms 2 focuses on confidence estimation and loss formulation.

Together, Algorithms 1 and 2 constitute the algorithmic backbone of our approach, jointly addressing the issue of how to dynamically identify and leverage sample-level dependency preferences among different modality combinations.

Algorithm 1 Feature Extraction and Modality Allocation

Require: Text X_t , Vision X_v , Audio X_a

Ensure: $\tilde{\mathbf{h}}_e$

Feature Extraction

- 1: $\mathbf{F}_t \leftarrow \text{BERT}(X_t)$
- 2: $\mathbf{F}_v \leftarrow \text{Transformer}(X_v)$
- 3: $\mathbf{F}_a \leftarrow \text{Transformer}(X_a)$

Hyperspherical Mapping

- 4: $\mathbf{h}_m \leftarrow \text{LayerNorm}(\mathbf{F}_m), m \in \{t, v, a\}$

Routing

- 5: $\mathbf{q}_{e,m} \leftarrow \text{softmax}(w_{e,m} \mathbf{h}_m), e \in \{|E|\}$
- 6: $\mathbf{R}_{e,m} \leftarrow \text{Top-k}(q_{e,m}, k)$

Modality Allocation

- 7: **for** Experts $e = 1$ to $|E|$ **do**
- 8: $\tilde{\mathbf{h}}_e \leftarrow \text{MLP}_e(\sum_m \mathbf{R}_{e,m} \cdot \mathbf{h}_m)$
- 9: **end for**

Algorithm 2 Reliability-Aware Expert Selection and Loss

Require: $\tilde{\mathbf{h}}_e$, Label y

Ensure: \hat{y} , \mathcal{L}_{att} , \mathcal{L}_{rep} , \mathcal{L}_{con}

Angular-Validation

```

1: if  $y \neq \text{None}$  then
2:    $\tilde{y} \leftarrow \text{softmax}(\text{GaussianSoftLabel}(y))$ 
3: end if
4: for  $e = 1$  to  $|E|$  do
5:    $\mu_c \leftarrow \mathcal{L}_{con}$ 
6:    $\cos \theta_{e,c} \leftarrow \mu_c^\top \text{LayerNorm}(\tilde{\mathbf{h}}_e)$ 
7:   if Training then
8:      $\mathbf{z}_{e,c} \leftarrow s \cdot \cos(\theta_{e,c} + \nu)$ 
9:   else
10:     $\mathbf{z}_{e,c} \leftarrow s \cdot \cos(\theta_{e,c})$ 
11:  end if
12:   $\mathbf{p}_e \leftarrow \text{softmax}(\mathbf{z}_{e,c})$ 
13: end for
14:  $\varpi_e \leftarrow \sum^C p_e(c) \cdot \tilde{y}$ 
15:  $K_e \leftarrow -\sum^C p_e(c) \log p_e(c)$ 

```

RAES

```

16:  $H_e \leftarrow \|\hat{\mathbf{h}}_e\|_2$ 
17:  $\pi_e \leftarrow H_e, K_e, \varpi_e$ 
18:  $\hat{h} \leftarrow \sum^{|E|} \text{softmax}(\pi_e) \tilde{\mathbf{h}}_e$ 
19:  $\hat{y} \leftarrow \text{MLP}(\hat{h})$ 

```

Loss

```

20:  $\mathcal{L}_{task} \leftarrow \text{MSE}(\hat{y}, y)$ 
21: for  $m, n \in \{t, v, a\}$ ,  $m \neq n$  do
22:    $A_{Ge}^{(m,n)} \leftarrow h_m h_n^\top$ 
23:    $A_{Gr}^{(m,n)} \leftarrow \mathcal{L}_{task}(m, n)$ 
24:    $A_{m,n} \leftarrow A_{Ge}^{(m,n)}, A_{Gr}^{(m,n)}$ 
25:    $\mathcal{L}_{att} \leftarrow A_{m,n} \text{KL}(\mathbf{q}_n, \mathbf{q}_m)$ 
26:    $\mathcal{L}_{rep} \leftarrow -(1 - A_{m,n}) \text{JS}(\mathbf{q}_n, \mathbf{q}_m)$ 
27: end for
28:  $\mathcal{L}_{con} \leftarrow \text{KL}(\mathbf{p}_e, \tilde{y})$ 

```
