

# Resource-Limited Joint Multimodal Sentiment Reasoning and Classification via Chain-of-Thought Enhancement and Distillation

Anonymous ACL submit

## Abstract

Current approaches for Multimodal Sentiment Analysis (MSA) primarily leverage the knowledge and reasoning capabilities of parameter-heavy (Multimodal) LLMs for classification, overlooking autonomous multimodal sentiment reasoning generation in resource-constrained environments. In this paper, we focus on the Resource-Limited Joint Multimodal Sentiment Reasoning and Classification task, JM-SRC, which simultaneously performs multimodal sentiment reasoning chain generation and sentiment classification only with a lightweight model. We propose a Multimodal Chain-of-Thought Reasoning Distillation model, MulCoT-RD, designed for JM-SRC that employs a "Teacher-Assistant-Student" distillation paradigm to address deployment constraints in resource-limited environments. We first leverage a high-performance Multimodal Large Language Model (MLLM) to generate the initial reasoning dataset and train a medium-sized assistant model with a multi-task learning mechanism. A lightweight student model is jointly trained to perform efficient multimodal sentiment reasoning generation and classification. Extensive experiments on four datasets demonstrate that MulCoT-RD<sup>1</sup> with only 3B parameters and achieves strong performance on JM-SRC, while exhibiting robust generalization and enhanced interpretability.

## 1 Introduction

With the proliferation of social media and multimedia content, Multimodal Sentiment Analysis (MSA) has emerged as a critical research area attracting significant academic and industry attention Yang et al. (2024); Amiriparian et al. (2024). MSA of text-image pairs can be categorized into coarse-grained and fine-grained approaches based on sentiment targets. Coarse-grained MSA (Yang et al.,

<sup>1</sup>The code and demo are available via <https://anonymous.4open.science/r/MulCoT-RD/>.

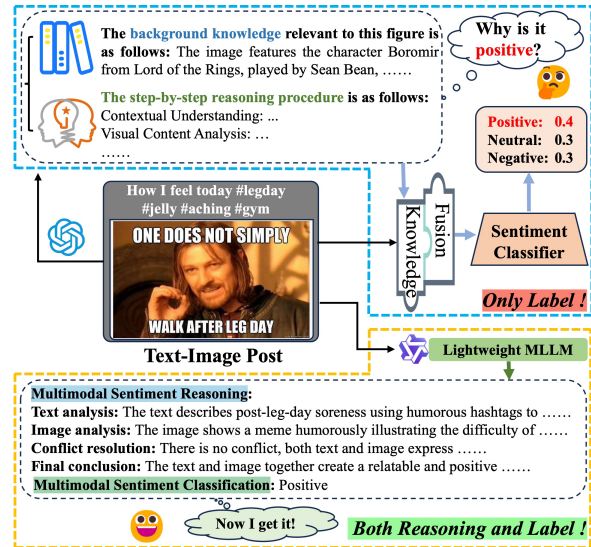


Figure 1: Leveraging reasoning (blue dashed line) vs. Generating reasoning chain (yellow dashed line).

2021; Zhang et al., 2023) identifies the overall sentiment of text-image pairs, while fine-grained MSA, or Multimodal Aspect-Based Sentiment Classification (MASC) (Zhou et al., 2023; Wang et al., 2024; Yang et al., 2025b), analyzes sentiment toward specific aspect terms within textual content.

Most existing methods enhance MSA through multimodal representation learning (Zhang et al., 2022; Manzoor et al., 2023) and fusion (Huang et al., 2020; Zhang et al., 2023), employing separate encoders to extract unimodal representations, then integrating them using fusion strategies such as gating mechanisms (Kumar and Vepa, 2020), cross-modal attention (Ju et al., 2021), and graph neural networks (Yang et al., 2021). While these approaches advance MSA performance, they face a fundamental limitation: inability to model intra-modal and cross-modal sentiment reasoning processes that explain why users experience particular sentiments. These models typically operate as "black boxes" for sentiment classification, obscuring the specific contributions of each modality and interaction mechanisms in sentiment decisions

064 due to their lack of explicit modeling of sentiment  
065 presentation and reasoning chain across modalities.

066 Building upon LLMs, Multimodal Large Lan-  
067 guage Models (MLLMs) (Hurst et al., 2024; Wu  
068 et al., 2024; Bai et al., 2025) demonstrate re-  
069 markable performance across diverse multimodal  
070 tasks, including recommendation systems (Ye et al.,  
071 2025), sentiment analysis (Wang et al., 2024), and  
072 mental health assessment (Zhang et al., 2024). As  
073 shown in Figure 1 (blue box), current methods  
074 leverage high-performing MLLMs, like GPT-4o,  
075 to inject world knowledge or Chain-of-Thought  
076 (CoT) (Wei et al., 2022) reasoning into pre-trained  
077 language models for MSA improvement (Wang  
078 et al., 2024; Li et al., 2025a), yet fail to trans-  
079 fer superior reasoning capabilities. Existing re-  
080 search (Li et al., 2025b) shows that lightweight  
081 MLLMs ( $\leq 3B$  parameters) exhibit limited CoT  
082 reasoning capabilities, necessitating reliance on  
083 models with superior reasoning abilities. How-  
084 ever, closed-source models incur substantial costs,  
085 while large-scale MLLMs require extensive compu-  
086 tational resources, limiting deployment in resource-  
087 constrained environments. Developing lightweight  
088 MLLMs (e.g., 3B parameters) that autonomously  
089 generate high-quality multimodal sentiment rea-  
090 soning while maintaining high MSA performance  
091 represents a major challenge, as highlighted in the  
092 yellow box of Figure 1.

093 To address these challenges, we focus on the  
094 **Resource-Limited Joint Multimodal Sentiment**  
095 **Reasoning and Classification (JMSRC)** task,  
096 which simultaneously performs multimodal senti-  
097 ment reasoning generation and classification using  
098 only a lightweight MLLM. We introduce the **Mul-**  
099 **timodal Chain-of-Thought Enhancement with**  
100 **Reasoning Distillation (MulCoT-RD)** framework  
101 for JMSRC, illustrated in Figure 2, while leverag-  
102 ing Reasoning Distillation (RD) with the Teacher-  
103 Assistant-Student pattern to enable lightweight  
104 MLLMs to autonomously generate high-quality  
105 sentiment reasoning (for the second challenge).  
106 The MulCoT-RD comprises two core modules. (1)  
107 **Multimodal CoT Enhancement Module:** We de-  
108 sign a two-stage module using structured prompt  
109 templates with task decomposition, reasoning guid-  
110 ance, conflict mediation steps, and adaptive retry  
111 control. It guides the high-performance closed-  
112 source or large-scale open-source MLLM as a  
113 teacher model to generate logically coherent multi-  
114 modal sentiment reasoning. (2) **Multimodal Senti-**  
115 **ment Reasoning Distillation Module:** Consider-

116 ing teacher model limitations in providing soft la-  
117 bels and intermediate representations, data scarcity,  
118 and inference costs, we introduce a medium-sized  
119 open-source MLLM as an assistant model, and use  
120 it to synthesize high-quality data. Through multi-  
121 task learning, the assistant model jointly enhances  
122 sentiment label prediction accuracy and reason-  
123 ing quality. For efficient deployment in resource-  
124 constrained environments, we employ joint opti-  
125 mization combining hard labels with soft labels  
126 from the assistant model to transfer reasoning capa-  
127 bilities to a lightweight student MLLM, achieving  
128 optimal balance among classification performance,  
129 interpretability, and deployment efficiency. Our  
130 contributions are summarized as follows:

- We focus on joint multimodal sentiment  
131 reasoning and classification in resource-  
132 constrained scenarios and construct a high-  
133 quality sentiment reasoning dataset. 134
- We propose the Multimodal Chain-of-  
135 Thought Enhancement with Reasoning  
136 Distillation, MulCoT-RD, framework for  
137 JMSRC. Multi-task learning and joint opti-  
138 mization improve the sentiment classification  
139 and reasoning capabilities of the model. 140
- Comprehensive experiments across multi-  
141 ple MSA datasets demonstrate that our  
142 lightweight 3B-parameter MLLM achieves  
143 superior sentiment classification performance  
144 while maintaining high interpretability. 145

## 2 Related Work 146

### 2.1 Multimodal Sentiment Analysis 147

148 The MSA development can be broadly divided into  
149 two stages: the era of pre-trained language mod-  
150 els (PLMs) and the era of large language models  
151 (LLMs). During the PLMs era, MSA methods typ-  
152 ically utilize a dedicated encoder for each modal-  
153 ity to extract representations, with a primary fo-  
154 cus on multimodal fusion and cross-modal align-  
155 ment. (Zhang et al., 2023; Xiao et al., 2023; Zhou  
156 et al., 2023). The emergence of LLMs has opened  
157 new possibilities for MSA. However, existing meth-  
158 ods typically rely on MLLMs to generate valuable  
159 knowledge (Wang et al., 2024) or reasoning (Pang  
160 et al., 2024; Li et al., 2025a), which is then injected  
161 into pre-trained language models to improve MSA,  
162 rather than enabling autonomous sentiment reason-  
163 ing. This results in limited interpretability. To  
164 our knowledge, Emotion-LLaMA is the first LLM-  
165 based model for multimodal emotion recognition

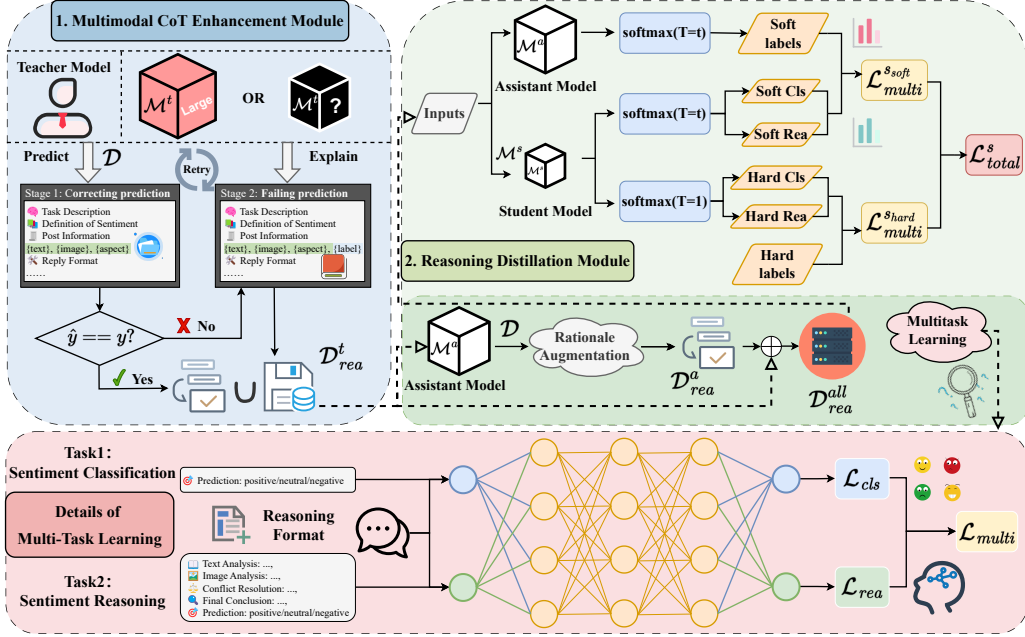


Figure 2: MulCoT-KD comprises two core modules, i.e., (1) Multimodal CoT Enhancement Module, (2) Reasoning Distillation Module (Assistant Model with Multi-Task Learning, Student Model with Joint Learning).

and explanation, but requires modality-specific representation learning, pre-training, and instruction tuning (Cheng et al., 2024). Models with superior reasoning capabilities are often computationally expensive or have large parameter counts that complicate deployment. We focus on using lightweight MLLM to simultaneously achieve efficient and autonomous generation of high-quality multimodal sentiment reasoning and classification.

## 2.2 Reasoning Distillation

Knowledge Distillation (KD) (Hinton et al., 2015) has proven effective for compressing language models by transferring predictive behaviors, such as soft labels or hidden representations, from larger teacher models to smaller student models. Current KD techniques for PLMs focus on distilling soft labels (Sanh et al., 2019; Gu et al., 2023) or representations (Wang et al., 2020b,a; Kim et al., 2022), but require access to the teacher model’s internal parameters. This dependency creates significant challenges when applying KD to closed-source LLMs. Reasoning distillation offers an alternative approach, enabling smaller student models to acquire reasoning capabilities by fine-tuning on reasoning processes from a teacher model instead of relying on soft labels (Magister et al., 2022; Li et al., 2023; Lee et al., 2024; Chenglin et al., 2024). In our work, we leverage an intermediate-sized model with multi-task learning as an assistant to both supplement soft-label distillation signals from the teacher model and generate higher-quality data

to address reasoning data scarcity.

## 3 Method

To achieve an effective integration of task performance, interpretability, and deployment efficiency, we introduce the Multimodal Chain-of-Thought Enhancement with Reasoning Distillation (MulCoT-RD) framework for JMSRC, as shown in Figure 2, comprising the Multimodal CoT Enhancement Module and the Reasoning Distillation Module.

### 3.1 Task Definition

Given a dataset  $\mathcal{D} = \{(x_i, L_i)\}_{i=1}^N$  containing  $N$  samples, each sample  $x_i$  consists of text  $T_i$ , image  $I_i$ , aspect term  $[A_i]$  (provided only in fine-grained MSA), and sentiment label  $L_i$ . The JMSRC task is formulated as follows:

$$\mathcal{M}(T_i, I_i, [A_i]) \Rightarrow (R_i, \hat{y}_i), \quad (1)$$

where  $R_i$  denotes the corresponding sentiment reasoning, and  $\hat{y}_i$  denotes the predicted sentiment label by MLLM  $\mathcal{M}$ .

### 3.2 Multimodal CoT Enhancement

We propose a two-stage multimodal CoT enhancement module to synthesize high-quality sentiment reasoning data. The corresponding prompts are illustrated in Figure 5 provided in Appendix F. **In the first stage**, we perform reasoning path generation in a label-free setting using a high-performance MLLM as the teacher model  $\mathcal{M}^t$ . We employ a structured CoT prompt template  $\mathcal{T}_{pre}$  for **prediction**, comprising the basic template  $\mathcal{T}_b$  (including

Task Description, Sentiment Definition, and Reasoning Format) and the specific prediction prompt  $\mathcal{P}_{pre}$ . This template guides the model through text analysis, image analysis, conflict resolution, and conclusion generation, ensuring logically coherent and interpretable reasoning.

$$c_i^{t_1}, \hat{y}_i^t = \mathcal{M}^t(x_i; \mathcal{T}_{pre}), \quad (2)$$

where  $c_i^{t_1}$  represents the CoT reasoning process generated in the first stage, and  $\hat{y}_i^t$  indicates the predicted sentiment label for the  $i$ -th sample.

For correctly predicted samples, the generated reasoning paths are directly retained for subsequent training, thereby constructing the first-stage training set,  $\mathcal{D}_{rea}^{t_1}$ .

$$\mathcal{D}_{rea}^{t_1} = \{(x_i, c_i^{t_1}, \hat{y}_i^t) \mid \hat{y}_i^t = L_i\}_{i=1}^{N_{t_1}}. \quad (3)$$

Misclassified samples often reflect complex cases with ambiguous boundaries or cross-modal conflicts, or semantic ambiguity. Guiding the model to learn causally consistent reasoning on these challenging examples can enhance its understanding and robustness in complex scenarios. Therefore, we design a second stage where, for samples with incorrect predictions, the ground truth label,  $L_i$ , is introduced and an explain template,  $\mathcal{T}_{exp}$ , is constructed to guide the model in generating a supervised reasoning process,  $c_i^{t_2}$ , conditioned on the correct label.

$$\begin{cases} c_i^{t_2}, \hat{y}_i^t = \mathcal{M}^t(x_i, L_i; \mathcal{T}_{exp}) \\ \mathcal{D}_{rea}^{t_2} = \{(x_i, c_i^{t_2}, L_i)\}_{i=1}^{N_{t_2}}, \end{cases} \quad (4)$$

where  $N = N_{t_1} + N_{t_2}$ ;  $\mathcal{T}_{exp}$  is constructed by the basic template,  $\mathcal{T}_b$ , and the specific reasoning prompt,  $\mathcal{P}_{exp}$ , as shown in Figure 5 provided in Appendix F.

The two-stage datasets are merged to obtain the reasoning dataset  $\mathcal{D}_{rea}^t = \mathcal{D}_{rea}^{t_1} \cup \mathcal{D}_{rea}^{t_2}$ . To improve sentiment reasoning and label prediction reliability, we introduce an adaptive replay controller, ARC, that automatically regenerates outputs when MLLMs produce incomplete structures or invalid labels until a valid result is obtained or the retry limit is reached, ensuring generation quality while controlling computational overhead. The details of ARC are provided in Appendix E.

### 3.3 Multimodal Sentiment Reasoning Distillation

Closed-source teacher models limit knowledge extraction due to restricted intermediate representations, while open-source models with strong reasoning often require large parameters (Li et al.,

2025b), hindering efficient deployment. To address multimodal sentiment reasoning data scarcity and the absence of soft labels, we introduce reasoning distillation (Lee et al., 2024) to train an **assistant model with multi-task learning**, as illustrated in the middle right of Figure 2, enhancing data diversity. A **student model with joint learning**, as shown in the upper right of Figure 2, adapts to resource-constrained environments while inheriting the assistant model’s sentiment reasoning and classification capabilities.

#### 3.3.1 Assistant Model with Multi-Task Learning

We propose a multi-task learning framework that shares hard parameters to train the assistant model,  $\mathcal{M}^a$ , for JMSRC that jointly optimizes two complementary tasks, including multimodal sentiment reasoning and classification, as shown in the lower part of Figure 2.

$$\mathcal{L} = \frac{-1}{B} \sum_{i=1}^B \sum_{j=1}^l \log P\left(y_j^{(i)} \mid y_{<j}^{(i)}, \mathcal{M}^a(x^{(i)})\right) \cdot I_{y_j^{(i)} \neq -100}, \quad (5)$$

where  $B$  denotes the batch size;  $l$  denotes the target sequence length of the  $i$ -th sample;  $P$  denotes the predicted probability of  $y_j^{(i)}$  at decoding step  $j$  based on  $y_{<j}^{(i)}$ ;  $I_{y_j^{(i)} \neq -100}$  indicates that only tokens whose labels are not equal to -100 (i.e., not masked) participate in the loss.

The overall loss function for training the assistant model is formulated as follows:

$$\mathcal{L}_{multi}^a = \lambda_{cls}^a \cdot \mathcal{L}_{cls}^a + \lambda_{rea}^a \cdot \mathcal{L}_{rea}^a, \quad (6)$$

where  $\lambda_{cls}^a$  and  $\lambda_{rea}^a$  are the weighting hyperparameters to ensure a balanced trade-off between two tasks. After training, we can obtain the trained assistant model,  $\overline{\mathcal{M}}^a$ .

Regarding data augmentation, given the limited capabilities of the assistant model, we only retain training samples for which sentiment can be correctly predicted through sentiment reasoning. See the material for more details.

$$\mathcal{D}_{rea}^a = \{(x_i, \hat{c}_i^a, \hat{y}_i^a) \mid \hat{y}_i^a = L_i\}_{i=1}^{N_a}, \quad (7)$$

where  $\hat{c}_i^a, \hat{y}_i^a = \overline{\mathcal{M}}^a(x_i; \mathcal{T}_{pre})$  and  $N_a < N$ .

The complete sentiment reasoning dataset is obtained, which is used to train a student model.

$$\mathcal{D}_{rea}^{all} = \mathcal{D}_{rea}^t \cup \mathcal{D}_{rea}^a. \quad (8)$$

### 3.3.2 Student Model with Joint Learning

To enable efficient deployment in resource-constrained environments, we employ a lightweight student MLLM,  $\mathcal{M}^s$ , trained through knowledge distillation. The student model jointly learns from two sources, including ground-truth labels (hard labels) for accurate prediction and probability distributions (soft labels) from the assistant model to capture its reasoning patterns. The dual supervision allows the student model to inherit the assistant model’s discriminative capabilities.

**Hard Label.** The student model undergoes fine-tuning using constructed reasoning data,  $\mathcal{D}_{rea}^{all}$ , enabling it to acquire step-by-step reasoning capabilities through reasoning distillation. The hard label loss is defined as follows:

$$\begin{cases} \mathcal{L}_{cls}^{shard} = \mathbb{E}_{\mathcal{D}_{rea}^{all}} \log P([x; L] | \mathcal{M}^s) \\ \mathcal{L}_{rea}^{shard} = \mathbb{E}_{\mathcal{D}_{rea}^{all}} \log P([x; c] | \mathcal{M}^s), \end{cases} \quad (9)$$

where  $P$  denotes the probability distribution;  $c$  represents the reasoning process. The losses  $\mathcal{L}_{cls}^{shard}$  and  $\mathcal{L}_{rea}^{shard}$  are used to train the student model to learn the direct mapping from multimodal input to sentiment labels and to generate coherent sentiment reasoning, respectively.

**Soft Label.** To address the black-box nature of closed-source MLLMs, the assistant model is employed as an intermediary to provide soft labels for distillation. Given an input  $x$ , the probability distribution  $p_k$  at the  $k$ -th position is obtained from the logit value  $z_k$  through a single forward pass followed by the softmax function. It is formally defined as:

$$p_k = \frac{\exp(z_k/\tau)}{\sum_j \exp(z_j/\tau)}, \quad (10)$$

where  $\tau$  denotes the temperature hyperparameter, which is used to control the smoothness of the distribution.

After obtaining the probability distributions  $p^a$  from  $\mathcal{M}^a$  and  $p^s$  from  $\mathcal{M}^s$ , we employ the Kullback–Leibler (KL) (Wu et al., 2025) divergence to minimize the discrepancy between the two distributions. It enables the student model to mimic the prediction behavior of the larger model. The training for soft label distillation is defined as follows:

$$\begin{cases} \mathcal{L}_{soft}(p^a, p^s) = \sum_k p_k^a \log \frac{p_k^a}{p_k^s} \\ \mathcal{L}_{cls}^{soft} = \mathcal{L}_{soft}(p_{cls}^a, p_{cls}^s) \\ \mathcal{L}_{rea}^{soft} = \mathcal{L}_{soft}(p_{rea}^a, p_{rea}^s). \end{cases} \quad (11)$$

**Joint Learning.** The student model training retains the multi-task learning. The overall hard-label loss and soft-label loss for the student model are defined as follows:

$$\begin{cases} \mathcal{L}_{multi}^{shard} = \lambda_{cls}^{shard} \cdot \mathcal{L}_{cls}^{shard} + \lambda_{rea}^{shard} \cdot \mathcal{L}_{rea}^{shard} \\ \mathcal{L}_{multi}^{soft} = \lambda_{cls}^{soft} \cdot \mathcal{L}_{cls}^{soft} + \lambda_{rea}^{soft} \cdot \mathcal{L}_{rea}^{soft} \end{cases} \quad (12)$$

where  $\lambda_{cls}^{shard}$ ,  $\lambda_{rea}^{shard}$ ,  $\lambda_{cls}^{soft}$ , and  $\lambda_{rea}^{soft}$  are hyperparameters that balance the contributions of classification loss and reasoning generation loss in the hard-label and soft-label multi-task learning objectives, respectively.

To jointly leverage hard-label and soft-label supervision, we define the total loss of the student model as follows.

$$\mathcal{L}_{total}^s = (1 - \lambda) \mathcal{L}_{multi}^{shard} + \lambda \mathcal{L}_{multi}^{soft}, \quad (13)$$

where  $\lambda$  is a hyperparameter that controls the balance between hard-label and soft-label supervision.

## 4 Experiments

### 4.1 Experimental Settings

#### 4.1.1 Datasets

We conduct experiments on both coarse-grained MSA, MVSA-Single and MVSA-Multiple datasets, preprocessed following (Liu et al., 2024) and fine-grained MSA, Twitter-2015 and Twitter-2017 datasets (Yu and Jiang, 2019). Table 1 presents the statistics of four datasets with the constructed sentiment reasoning data for JMSRC.

Dataset	Train	Dev	Test	Train <sup>g+</sup>	Train <sup>q+</sup>
MVSA-Single	3608	451	452	6483	6350
MVSA-Multiple	13619	1702	1702	23424	23697
Twitter-2015	3179	1122	1037	6166	6218
Twitter-2017	3562	1176	1234	6652	6871

Table 1: Statistics of datasets.  $g+$  and  $q+$  represent the teacher models GPT-4o-mini (Hurst et al., 2024) and Qwen2.5-VL-72B (Bai et al., 2025), respectively.

#### 4.1.2 Model Selection

To build an efficient hierarchical reasoning distillation, we design four distillation architectures, as summarized in Table 2. Note that, while our model selection is limited, experimental results clearly demonstrate the effectiveness of MulCoT-RD. See the Appendix B for more details.

#### 4.1.3 Implementation Details

We train our models on NVIDIA RTX A6000 GPUs using the AdamW optimizer (Loshchilov

ID	Teacher Model	Assistant Model	Student Model
1	<b>GPT-4o-mini</b>	Qwen3-VL-8B	Qwen3-VL-2B
2		Qwen2.5-VL-7B	Qwen2.5-VL-3B
3	<b>Qwen2.5-VL-72B</b>	Qwen3-VL-8B	Qwen3-VL-2B
4		Qwen2.5-VL-7B	Qwen2.5-VL-3B

Table 2: Four reasoning distillation architectures.

and Hutter, 2017). During training, we set the initial learning rate to  $3e-4$  and employ a dynamic adjustment strategy: if the validation set performance does not improve for two consecutive epochs, we halve the learning rate until it reaches a minimum of  $1e-6$ . Due to resource limitations, we set the batch size to 2 and train for a maximum of 20 epochs. To mitigate instability caused by small batch sizes, we use gradient accumulation, updating parameters every 20 steps. The multi-task learning hyperparameters  $\lambda_{rea}^a$ ,  $\lambda_{rea}^{shard}$ ,  $\lambda_{rea}^{soft}$  and  $\lambda_{cls}^a$ ,  $\lambda_{cls}^{shard}$ ,  $\lambda_{cls}^{soft}$  are set to 0.8 and 0.2, respectively, while the knowledge distillation coefficient  $\lambda$  is set to 0.3. Detailed explanations and configurations can be found in Appendix D

#### 4.1.4 Evaluation Metrics

In line with previous work (Chen et al., 2024), we evaluate model performance of classification on coarse-grained MSA using Accuracy (**Acc**) and Weighted F1 (**w-F1**). For fine-grained MSA (MASC), we follow previous studies (Zhou et al., 2023) and adopt Accuracy and Macro F1 (**m-F1**) as evaluation metrics. For the sentiment reasoning task, we employ comprehensive metrics including sentence embedding-based cosine similarity (**Sim**) (Reimers and Gurevych, 2019), **METEOR** (Banerjee and Lavie, 2005), **BLEU** (Papineni et al., 2002), **ROUGE-L** (Lin, 2004), and Distinct-N1/N2 (**Dist-1/2**) (Li et al., 2015).

## 4.2 Baselines

We compare popular models on **coarse-grained MSA** with MulCoT-RD, including **MultiSentiNet** (Xu and Mao, 2017), **HSAN** (Xu, 2017), **CoMN-Hop6** (Xu et al., 2018), **MGNNS** (Yang et al., 2021), **CLMLF** (Li et al., 2022), **MVCN** (Wei et al., 2023), **D<sup>2</sup>R** (Chen et al., 2024). For **fine-grained MSA**, involving **ESAFN** (Yu et al., 2019), **TomBERT** (Yu and Jiang, 2019), **CapTrBERT** (Khan and Fu, 2021), **JML** (Ju et al., 2021), **VLP-MABSA** (Ling et al., 2022), **CMMT** (Yang et al., 2022), **AoM** (Zhou et al., 2023), **AETS** (Zhu et al., 2025). **Emotion-LLaMA** (Cheng et al., 2024) employs pretraining and instruction tuning based on LLaMA2-7B-Chat to enhance multimodal

emotion recognition and explanation. **Qwen3-VL-8B-Thinking** (Yang et al., 2025a) features Interleaved-MRoPE and DeepStack for powerful spatial-temporal reasoning. Detailed descriptions can be found in Appendix C.

## 4.3 Main Results

Unlike previous models that only perform multimodal sentiment classification, our model enables joint sentiment reasoning and classification. We conduct experiments on both multimodal sentiment classification and reasoning tasks.

### 4.3.1 Results of Multimodal Sentiment Classification

**Performance on coarse-grained MSA.** Table 3 presents the comparison results on the coarse-grained MSA task. MulCoT-RD outperforms both the second-best model (Emotion-LLaMA) and the previous state-of-the-art model (D<sup>2</sup>R) on the MVSA-Single and MVSA-Multiple datasets, achieving substantial improvements. It highlights the benefits of explicitly modeling intra-modal sentiment structures and cross-modal reasoning processes. Notably, although the teacher model has greater parameter capacity, its lack of task-specific fine-tuning for MSA leads to suboptimal modeling of cross-modal emotional relations, making it inferior to the assistant model optimized with task-oriented objectives. Moreover, the student model outperforms the assistant model in certain cases, likely due to benefiting from the augmented training data generated by the assistant, which improves its generalization and robustness.

**Performance on MASC.** As shown in Table 4, the MulCoT-RD(asst) model (with Qwen2.5-VL-72B as the teacher) achieves the best overall performance. Compared to the second-best models AoM and AETS, MulCoT-RD(asst) exhibits a slight decrease in accuracy on the Twitter-2017 dataset by 1.4% and 1.6%, respectively, but consistently achieves the highest scores across all other evaluation metrics. We attribute this to two primary reasons. First, the Twitter-2017 dataset contains a large number of unparseable and unrecognizable symbols (Peng et al., 2024), including emojis that are commonly used on Twitter. These symbols may mislead the model by obscuring emotional semantics during reasoning, thereby slightly reducing accuracy. Second, MulCoT-RD(asst) is fine-tuned using LoRA, whereas most existing SOTA methods, such as AoM and AETS, adopt full-parameter fine-tuning. This limits the extent of parameter up-

Model	Venue	MVSA-S		MVSA-M	
		Acc	w-F1	Acc	w-F1
MultiSentinet	CIKM'17	69.8	69.8	68.9	68.1
HSAN	ISI'17	69.9	66.9	68.0	67.8
CoMN-Hop6	SIGIR'18	70.5	70.0	68.9	68.8
MGNNS	ACL'21	73.8	72.7	72.5	69.3
CLMLF	NAACL'21	75.3	73.5	72.0	69.8
MVCN	ACL'23	76.1	74.6	72.1	70.0
D <sup>2</sup> R	EMNLP'24	76.7	75.6	71.6	70.9
Emotion-LLaMA <sup>†</sup>	NeurIPS'24	82.7	81.8	75.6	<b>75.2</b>
Open-Flamingo <sup>‡</sup>	ICML'25	66.3	-	68.7	-
Qwen2.5-VL-3B*	Student	62.8	66.4	74.2	70.7
Qwen2.5-VL-7B*	Assistant	67.7	69.6	74.7	70.9
3-VL-8B-Thinking*		72.3	71.5	70.2	69.4
GPT-4o-mini*	Teacher <sup>1</sup>	76.7	75.6	71.6	71.4
<b>MulCoT-RD(asst)</b>		<b>83.6</b>	82.8	75.7	72.9
<b>MulCoT-RD(stu)</b>		82.7	82.3	<u>76.9</u>	74.2
Qwen2.5-VL-72B*	Teacher <sup>2</sup>	67.9	70.8	74.2	71.8
<b>MulCoT-RD(asst)</b>		83.2	82.1	<u>76.9</u>	73.8
<b>MulCoT-RD(stu)</b>		<u>83.4</u>	<b>83.2</b>	<b>77.2</b>	<u>74.4</u>

Table 3: Results for coarse-grained MSA. Models above the middle line are small models fully fine-tuned, while those below are (M)LLMs fine-tuned with LoRA. <sup>†</sup> denotes the results reproduced by us using models re-trained on our datasets. <sup>‡</sup> indicates the 16-shot performance under In-Context Learning (ICL). The best results are bold-typed and the second best ones are underlined. \* means the zero-shot performance.

dates during task adaptation, resulting in smaller performance gains compared to full fine-tuning (Biderman et al., 2024). Given this, we believe our proposed method remains effective for MASC.

Notably, the student model of MulCoT-RD contains only 3B parameters, significantly fewer than the large multimodal architecture of Emotion-LLaMA (Cheng et al., 2024), which combines LLaMA2-7B-chat with encoders like EVA, CLIP, VideoMAE, and HuBERT-large. Despite its smaller size, MulCoT-RD(stu) outperforms Emotion-LLaMA on multiple benchmarks, demonstrating superior efficiency and strong applicability in resource-constrained settings.

### 4.3.2 Evaluation of Sentiment Reasoning

MulCoT-RD achieves efficient and effective sentiment reasoning. We evaluate the reasoning performance of the student and assistant models, as well as Emotion-LLaMA, using the sentiment reasoning process from the teacher model as gold-standard references (exemplified by GPT-4o-mini), with results presented in Table 5. Our models achieve a comprehensive performance advantage over Emotion-LLaMA across all key reasoning metrics. The results demonstrate high-quality sentiment reasoning generation across multiple evaluation metrics. Cosine similarity (Sim) consis-

Model	Venue	Twitter-15		Twitter-17	
		Acc	m-F1	Acc	m-F1
ESAFN	TASLP'20	73.4	67.4	67.8	64.2
TomBERT	IJCAI'19	77.2	71.8	70.5	68.0
CapTrBERT	ACM MM'21	78.0	73.2	72.3	70.2
JML	EMNLP'21	78.7	-	72.7	-
VLP-MABSA	ACL'22	78.6	73.8	73.8	71.8
CMMT	IPM'22	77.9	-	73.8	-
AoM	ACL'23	80.2	<u>75.9</u>	<u>76.4</u>	<u>75.0</u>
AETS	AAAI'25	79.5	-	<b>76.6</b>	-
Emotion-LLaMA <sup>†</sup>	NeurIPS'24	73.9	70.2	69.2	67.9
Open-Flamingo <sup>‡</sup>	ICML'25	70.4	-	62.6	-
Qwen2.5-VL-3B*	Student	48.9	49.7	56.8	55.6
Qwen2.5-VL-7B*	Assistant	58.3	55.6	58.6	57.6
3-VL-8B-Thinking*		59.2	55.8	60.4	58.9
GPT-4o-mini*	Teacher <sup>1</sup>	49.4	37.6	54.0	52.8
<b>MulCoT-RD(asst)</b>		<u>80.7</u>	75.3	74.6	74.6
<b>MulCoT-RD(stu)</b>		80.4	75.2	74.0	73.3
Qwen2.5-VL-72B*	Teacher <sup>2</sup>	59.5	57.1	63.9	63.4
<b>MulCoT-RD(asst)</b>		<b>80.8</b>	<b>77.2</b>	75.0	<b>75.1</b>
<b>MulCoT-RD(stu)</b>		80.5	75.1	74.3	74.1

Table 4: Results of different methods for MASC. “-” means it does not exist in the original paper.

tently exceeds 90% across all models, confirming strong semantic alignment between generated and gold-standard reasoning chains. METEOR scores ranging from 45.4% to 59.8% further indicate substantial paraphrase-level and lexical overlap. While BLEU and ROUGE-L show some fluctuations, coarse-grained MSA variants generally outperform fine-grained MSA, reflecting better surface-form alignment. Distinct-N1 and Distinct-N2 scores remain approximately 49% and 80%, respectively, indicating that the generated reasoning maintains high linguistic diversity, enhancing the interpretability and robustness of reasoning tasks. Case studies can be found in the Appendix I.

Model	Dataset	Sim	Meteor	Bleu	Rouge-L	Dist-1	Dist-2
ELLA	MVSA-S	87.6	35.9	14.6	35.1	49.8	80.2
	MVSA-M	84.7	36.0	15.9	35.9	52.5	83.7
	Twitter-15	86.3	38.6	18.3	39.3	42.7	72.9
	Twitter-17	86.6	38.1	17.6	38.2	43.0	73.1
Asst	MVSA-S	92.6	59.8	47.8	55.0	49.8	80.2
	MVSA-M	93.0	57.4	48.1	57.2	48.6	79.4
	Twitter-15	92.9	54.6	43.0	58.3	42.4	72.9
	Twitter-17	90.5	51.2	35.9	53.3	45.2	74.1
Stu	MVSA-S	92.2	47.3	58.8	54.2	49.8	80.2
	MVSA-M	92.1	56.8	46.7	55.8	49.5	80.3
	Twitter-15	90.3	45.4	28.2	46.0	49.5	79.9
	Twitter-17	90.0	49.2	33.1	50.8	45.2	74.1

Table 5: Evaluation results of generated reasoning from Emotion-LLaMA, assistant and student models.

## 4.4 Ablation Study

In this section, we investigate the impact of each MulCoT-RD component, with results presented in

Table 6. When we only use the text modality (**w/o Img**), the model performs worse on all metrics compared to the complete model, highlighting the importance of incorporating visual modality. Similarly, when we remove the text modality (**w/o Text**), the model has a significant performance drop on all datasets. The decline, more severe than w/o Img, highlights the key role of text and the necessity of multimodal integration. **w/o Rea** means to remove the multi-task learning paradigm and exclude the sentiment reasoning task from the training process, leading to a general performance drop. It highlights the importance of deeply modeling intra-modal and cross-modal sentiment reasoning. Note that all the above ablation experiments are conducted on the assistant model. **w/o Asst** omits the assistant model, removing the use of soft labels in the distillation process and reducing the scale and diversity of training data. This leads to a notable performance drop across all datasets, demonstrating the effectiveness of the teacher–assistant–student hierarchical distillation framework for JMSRC.

Method	MVSA-S		MVSA-M		Twitter-15		Twitter-17	
	Acc	w-F1	Acc	w-F1	Acc	m-F1	Acc	w-F1
MulCoT-RD	<b>83.2</b>	<b>82.1</b>	<b>76.9</b>	73.8	<b>80.8</b>	<b>77.2</b>	<b>75.0</b>	<b>75.1</b>
w/o Img	79.4	77.7	73.7	73.0	78.4	72.5	73.5	73.5
w/o Txt	77.9	77.1	66.2	67.7	65.6	56.6	64.6	59.4
w/o CoT	79.9	79.7	74.2	73.1	79.9	75.5	74.2	73.4
w/o Asst	81.9	81.3	75.2	<b>74.1</b>	79.3	72.3	73.7	73.3

Table 6: The performance comparison of our full model and its ablated methods under the setting where Qwen2.5-VL-72B serves as the teacher model.

#### 4.5 Efficiency of MulCoT-RD

To further demonstrate the practicability of MulCoT-RD, we provide the model efficiency comparison in Figure 3. We find that, on the MVSA-Single and Twitter-2015 datasets, our distilled student models (Qwen2.5-VL-3B and Qwen3-VL-2B) achieve significantly lower inference latency and GPU memory usage than Emotion-LLaMA, while obtaining notable improvements in Accuracy. Specifically, in Twitter-2015, Qwen3-VL-2B achieves the 8.79% accuracy improvement (80.4 vs. 73.9) with 5.81x fewer parameters and 0.33x faster inference speed compared to Emotion-LLaMA.

#### 4.6 Robustness of MulCoT-RD

To validate the robustness of our approach across different backbones, we conducted the base-model adaptation study by replacing the Qwen2.5-VL series with the Qwen3-VL series (as shown in Ta-

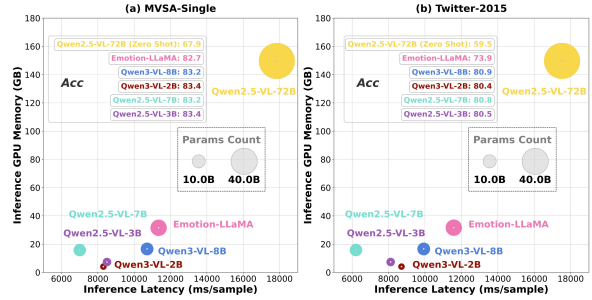


Figure 3: Efficiency comparison on MVSA-Single and Twitter-2015. Metrics are measured on the batch size as 1 and all samples are from the test set. Note that the models are trained under the paradigm where Qwen2.5-VL-72B serves as the teacher model. Detailed efficiency comparison results are provided in Appendix H.

ble 7) and the Flan-T5 series (see Appendix G). The results demonstrate that the models maintain strong performance across different backbones, illustrating the robustness and adaptability of MulCoT-RD.

Qwen3-VL	MVSA-S	MVSA-M	Twitter-15	Twitter-17
	Acc	w-F1	Acc	m-F1
<b>8B (asst)</b> <sup>1</sup>	<u>83.6</u>	<u>82.9</u>	76.0	73.5
<b>2B (stu)</b> <sup>1</sup>	<b>84.1</b>	<b>83.7</b>	<u>76.7</u>	<u>74.1</u>
<b>8B (asst)</b> <sup>2</sup>	83.2	82.6	<u>76.7</u>	73.2
<b>2B (stu)</b> <sup>2</sup>	83.4	82.4	<b>76.8</b>	<b>74.2</b>

Table 7: Performance of Qwen3-VL-based models on coarse-grained MSA and MASC. The best results are bold-typed, and the second best ones are underlined. 1 and 2 indicate the Teacher models, with 1 being GPT-4o-mini and 2 being Qwen2.5-VL-72B.

## 5 Conclusion

We focus on Joint Multimodal Sentiment Reasoning and Classification, JMSRC, in the resource-limited scenario that simultaneously generates multimodal reasoning chains and sentiment predictions. To address the dual challenges of reasoning interpretability and efficient deployment, we introduce MulCoT-RD, a unified framework combining structured CoT enhancement with reasoning distillation. Through a hierarchical teacher-assistant-student paradigm and joint multi-task learning, our method enables lightweight models to autonomously perform high-quality sentiment reasoning and classification. Extensive experiments across four datasets demonstrate the effectiveness and robustness of MulCoT-RD. In future work, we plan to incorporate direct preference optimization (DPO) with high- and low-quality reasoning sample filtering to further enhance the model’s emotional reasoning quality and classification performance.

## 600 Limitations

601 MulCoT-RD demonstrates strong performance on  
602 joint multimodal sentiment reasoning and clas-  
603 sification with a lightweight model, yet several  
604 limitations exist. Our approach depends on high-  
605 performance teacher models (e.g., GPT-4o-mini) to  
606 synthesize the initial reasoning dataset via carefully  
607 crafted Chain-of-Thought prompts. This genera-  
608 tion process entails considerable inference costs  
609 and is sensitive to the consistency and quality of  
610 the teacher model’s outputs. The evaluation of  
611 generated reasoning chains relies mainly on auto-  
612 matic metrics, including n-gram overlap measures  
613 and embedding-based cosine similarity. These  
614 metrics may fail to adequately assess semantic  
615 coherence, factual accuracy, logical consistency,  
616 or human-perceived interpretability. Due to re-  
617 source constraints, no large-scale human evalua-  
618 tion of reasoning quality was conducted. Although  
619 the student model is lightweight and targeted at  
620 resource-constrained deployment, it is built on  
621 a vision-language architecture that requires mul-  
622 timodal processing capabilities. This results in  
623 higher memory and computational demands than  
624 unimodal text-only models, which may limit de-  
625 ployment on extremely low-resource devices with  
626 minimal hardware acceleration support. We leave  
627 the exploration of more cost-effective reasoning  
628 data generation, comprehensive human evaluation  
629 of reasoning quality, multilingual generalization,  
630 and deployment on edge devices for future work.

## 631 Ethical considerations

632 This work constructs a multimodal sentiment rea-  
633 soning dataset by augmenting existing publicly  
634 available multimodal sentiment analysis datasets  
635 with model-generated reasoning produced by  
636 closed-source large multimodal language models  
637 (MLLMs). All source datasets are publicly avail-  
638 able and are used in accordance with their origi-  
639 nal licenses and terms of use. The closed-source  
640 MLLMs are accessed exclusively through offi-  
641 cial APIs, and no private, proprietary, or user-  
642 identifiable data are accessed or collected.

643 The proposed dataset does not introduce new per-  
644 sonal information beyond what is already present  
645 in the original public datasets, nor does it attempt  
646 to infer sensitive personal attributes. We acknowl-  
647 edge that sentiment interpretation is inherently sub-  
648 jective and that model-generated reasoning may  
649 reflect biases present in the underlying models. To

650 reduce this risk, we employ structured prompting  
651 and verification procedures to improve annotation  
652 consistency and reliability.

653 The dataset is intended solely for research pur-  
654 poses in multimodal sentiment understanding and  
655 reasoning. We believe that this work adheres to  
656 established ethical standards for data usage and re-  
657 sponsible application of large-scale models, and  
658 we do not anticipate foreseeable misuse beyond the  
659 scope of existing sentiment analysis research.

## References 660

- 661 Shahin Amiriparian, Lukas Christ, Alexander Kathan,  
662 Maurice Gerczuk, Niklas Müller, Steffen Klug,  
663 Lukas Stappen, Andreas König, Erik Cambria,  
664 Björn W Schuller, and 1 others. 2024. The muse  
665 2024 multimodal sentiment analysis challenge: So-  
666 cial perception and humor recognition. In *Proceed-  
667 ings of the 5th on Multimodal Sentiment Analysis  
668 Challenge and Workshop: Social Perception and Hu-  
669 mor*, pages 1–9.
- 670 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
671 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie  
672 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl  
673 technical report. *arXiv preprint arXiv:2502.13923*.
- 674 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An  
675 automatic metric for mt evaluation with improved re-  
676 lation with human judgments. In *Proceedings of  
677 the acl workshop on intrinsic and extrinsic evaluation  
678 measures for machine translation and/or summariza-  
679 tion*, pages 65–72.
- 680 Dan Biderman, Jacob Portes, Jose Javier Gonzalez Ortiz,  
681 Mansheej Paul, Philip Greengard, Connor Jennings,  
682 Daniel King, Sam Havens, Vitaliy Chiley, Jonathan  
683 Frankle, and 1 others. 2024. Lora learns less and  
684 forgets less. *arXiv preprint arXiv:2405.09673*.
- 685 Yifan Chen, Kuntao Li, Weixing Mai, Qiaofeng Wu,  
686 Yun Xue, and Fenghuan Li. 2024. D2r: Dual-branch  
687 dynamic routing network for multimodal sentiment  
688 detection. In *Proceedings of the 2024 Conference on  
689 Empirical Methods in Natural Language Processing*,  
690 pages 3536–3547.
- 691 Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Kai Wang,  
692 Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and  
693 Alexander Hauptmann. 2024. Emotion-llama: Mul-  
694 timodal emotion recognition and reasoning with in-  
695 struction tuning. *Advances in Neural Information  
696 Processing Systems*, 37:110805–110853.
- 697 Li Chenglin, Qianglong Chen, Liangyue Li, Caiyu  
698 Wang, Feng Tao, Yicheng Li, Zulong Chen, and Yin  
699 Zhang. 2024. Mixed distillation helps smaller lan-  
700 guage models reason better. In *Findings of the Asso-  
701 ciation for Computational Linguistics: EMNLP 2024*,  
702 pages 1673–1690.

703	Yanqi Dai, Zebin You, Dong Jing, Yutian Luo, Nanyi Fei, Guoxing Yang, and Zhiwu Lu. 2024. Cotbal: Comprehensive task balancing for multi-task visual instruction tuning. <i>arXiv preprint arXiv:2403.04343</i> .	756
704		757
705		758
706		759
707	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023. Minillm: Knowledge distillation of large language models. <i>arXiv preprint arXiv:2306.08543</i> .	760
708		761
709		762
710	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	763
711		764
712		765
713	Jian Huang, Jianhua Tao, Bin Liu, Zheng Lian, and Mingyue Niu. 2020. Multimodal transformer fusion for continuous emotion recognition. In <i>ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 3507–3511. IEEE.	766
714		767
715		768
716		769
717		770
718		771
719	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .	772
720		773
721		774
722		775
723		776
724	Xincheng Ju, Dong Zhang, Rong Xiao, Junhui Li, Shoushan Li, Min Zhang, and Guodong Zhou. 2021. Joint multi-modal aspect-sentiment analysis with auxiliary cross-modal relation detection. In <i>Proceedings of the 2021 conference on empirical methods in natural language processing</i> , pages 4395–4405.	777
725		778
726		779
727		780
728		781
729		782
730	Zaid Khan and Yun Fu. 2021. Exploiting bert for multimodal target sentiment classification through input space translation. In <i>Proceedings of the 29th ACM international conference on multimedia</i> , pages 3034–3042.	783
731		784
732		785
733		786
734		787
735	Junho Kim, Jun-Hyung Park, Mingyu Lee, Wing-Lam Mok, Joon-Young Choi, and SangKeun Lee. 2022. Tutoring helps students learn better: Improving knowledge distillation for bert with tutor network. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 7371–7382.	788
736		789
737		790
738		791
739		792
740		793
741		794
742	Ayush Kumar and Jithendra Vepa. 2020. Gated mechanism for attention based multi modal sentiment analysis. In <i>ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 4477–4481. IEEE.	795
743		796
744		797
745		798
746		799
747	Hojae Lee, Junho Kim, and SangKeun Lee. 2024. Mentor-kd: Making small language models better multi-step reasoners. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 17643–17658.	800
748		801
749		802
750		803
751		804
752	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. <i>arXiv preprint arXiv:1510.03055</i> .	805
753		806
754		807
755		808
		809
		810
		811
	Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023. Symbolic chain-of-thought distillation: Small models can also "think" step-by-step. <i>arXiv preprint arXiv:2306.14050</i> .	
	Yan Li, Xiangyuan Lan, Haifeng Chen, Ke Lu, and Dongmei Jiang. 2025a. Multimodal pear chain-of-thought reasoning for multimodal sentiment analysis. <i>ACM Transactions on Multimedia Computing, Communications and Applications</i> , 20(9):1–23.	
	Yuetai Li, Xiang Yue, Zhangchen Xu, Fengqing Jiang, Luyao Niu, Bill Yuchen Lin, Bhaskar Ramasubramanian, and Radha Poovendran. 2025b. Small models struggle to learn from strong reasoners. <i>arXiv preprint arXiv:2502.12143</i> .	
	Zhen Li, Bing Xu, Conghui Zhu, and Tiejun Zhao. 2022. Clmlf: A contrastive learning and multi-layer fusion method for multimodal sentiment detection. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 2282–2294.	
	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	
	Yan Ling, Jianfei Yu, and Rui Xia. 2022. Vision-language pre-training for multimodal aspect-based sentiment analysis. <i>arXiv preprint arXiv:2204.07955</i> .	
	Wuchao Liu, Wengen Li, Yu-Ping Ruan, Yulou Shu, Juntao Chen, Yina Li, Caili Yu, Yichao Zhang, Jihong Guan, and Shuigeng Zhou. 2024. Weakly correlated multimodal sentiment analysis: New dataset and topic-oriented model. <i>IEEE Transactions on Affective Computing</i> , 15(4):2070–2082.	
	Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. <i>arXiv preprint arXiv:1711.05101</i> .	
	Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn. 2022. Teaching small language models to reason. <i>arXiv preprint arXiv:2212.08410</i> .	
	Muhammad Arslan Manzoor, Sarah Albarri, Ziting Xian, Zaiqiao Meng, Preslav Nakov, and Shangsong Liang. 2023. Multimodality representation learning: A survey on evolution, pretraining and its applications. <i>ACM Transactions on Multimedia Computing, Communications and Applications</i> , 20(3):1–34.	
	Ning Pang, Wansen Wu, Yue Hu, Kai Xu, Qunjun Yin, and Long Qin. 2024. Enhancing multimodal sentiment analysis via learning from large language model. In <i>2024 IEEE International Conference on Multimedia and Expo (ICME)</i> , pages 1–6. IEEE.	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	



language models for multimodal sequential recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 13069–13077.

Jianfei Yu and Jing Jiang. 2019. Adapting bert for target-oriented multimodal sentiment classification. *IJCAI*.

Jianfei Yu, Jing Jiang, and Rui Xia. 2019. Entity-sensitive attention and fusion network for entity-level multimodal sentiment classification. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:429–439.

Haoyu Zhang, Yu Wang, Guanghao Yin, Kejun Liu, Yuanyuan Liu, and Tianshu Yu. 2023. Learning language-guided adaptive hyper-modality representation for multimodal sentiment analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 756–767.

Yazhou Zhang, Prayag Tiwari, Lu Rong, Rui Chen, Nojoom A AlNajem, and M Shamim Hossain. 2022. Affective interaction: Attentive representation learning for multi-modal sentiment classification. *ACM Transactions on Multimedia Computing, Communications and Applications*, 18(3s):1–23.

Yiqun Zhang, Xiaocui Yang, Xiaobai Li, Siyuan Yu, Yi Luan, Shi Feng, Daling Wang, and Yifei Zhang. 2024. Psydraw: A multi-agent multimodal system for mental health screening in left-behind children. *arXiv preprint arXiv:2412.14769*.

Ru Zhou, Wenya Guo, Xumeng Liu, Shenglong Yu, Ying Zhang, and Xiaojie Yuan. 2023. Aom: Detecting aspect-oriented information for multimodal aspect-based sentiment analysis. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8184–8196.

Linlin Zhu, Heli Sun, Qunshu Gao, Yuze Liu, and Liang He. 2025. Aspect enhancement and text simplification in multimodal aspect-based sentiment analysis for multi-aspect and multi-sentiment scenarios. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 1683–1691.

## A Data expansion with Assistant Model

After training the assistant model, we apply it to perform inference on the **original training set only**, explicitly excluding the validation and test sets to prevent any risk of label leakage. During this process, we retain only those samples whose predicted sentiment labels match the ground truth. These correctly predicted samples are then merged with the original training set to construct an expanded dataset, which is subsequently used for training the student model. Detailed results of the data expansion are presented in Table 8.

Dataset	Samples	GPT-4o-mini			Qwen2.5-VL-72B		
		Acc	w-F1	m-F1	Acc	w-F1	m-F1
MVSA-S	3608	79.7	79.7	69.9	76.0	77.1	66.9
MVSA-M	13619	72.0	68.0	55.3	74.0	70.5	60.6
Twitter-15	3179	94.0	94.1	92.6	95.6	95.6	94.6
Twitter-17	3562	86.8	86.7	86.4	92.9	92.9	93.4

Table 8: Performance of the Assistant Model (Qwen2.5-VL-7B) on Training Sets During Data Expansion, Guided by Different Teacher Models.

This strategy significantly increases the scale and diversity of the training data, broadens the coverage of sentiment label distributions, and incurs no additional manual annotation cost. It equips the student model with richer and higher-quality learning signals, effectively mitigating the challenge of limited annotated data commonly encountered in multimodal sentiment analysis tasks.

## B Model Selection

To construct a hierarchical reasoning distillation framework for achieving efficient joint multimodal sentiment reasoning and classification (JMSRC), we carefully select the following models as the teacher model, the assistant model, and the student model. Table 9 shows the specific model selections and their characteristics.

Role	Model	Access	Release Date
Teacher	GPT-4o-mini	Closed	2024.07
	Qwen2.5-VL-72B	Open	2025.02
Assistant	Qwen3-VL-8B	Open	2025.10
	Qwen2.5-VL-7B	Open	2025.02
Student	Qwen3-VL-2B	Open	2025.10
	Qwen2.5-VL-3B	Open	2025.02

Table 9: Model Selection and Characteristics.

## C Baselines

**Methods for coarse-grained MSA.** 1) **MultiSen-tiNet** (Xu and Mao, 2017) is a deep attention-based semantic network for multimodal sentiment analysis. 2) **HSAN** (Xu, 2017) is a hierarchical semantic attentional network based on image captions for multimodal sentiment analysis. 3) **CoMN-Hop6** (Xu et al., 2018) utilizes co-memory network to iteratively model the interactions between multiple modalities. 4) **MGNNS** (Yang et al., 2021) adopts multi-channel graph neural networks with sentiment-awareness for image-text sentiment detection. 5) **CLMLF** (Li et al., 2022) proposes a contrastive learning and multi-layer fusion method

for multimodal sentiment detection. 6) **MVCN** (Wei et al., 2023) designs a multi-view calibration network to solve the modality heterogeneity for multimodal sentiment detection. 7) **D<sup>2</sup>R** (Chen et al., 2024) proposes a dual-branch dynamic routing network to enhance multimodal sentiment detection by effectively modeling cross-modal interactions. 8) **Emotion-LLaMA** (Cheng et al., 2024) employs a specialized emotion tokenizer and instruction fine-tuning based on the LLaMA2-7B-chat to enhance multimodal emotion recognition. 9) **Qwen3-VL-8B-Thinking** (Yang et al., 2025a) features Interleaved-MRoPE and DeepStack for powerful spatial-temporal reasoning, plus precise Text–Timestamp Alignment. With native 256K context, it excels at complex STEM and logic tasks.

**Methods for fine-grained MSA.** 1) **ESAFN** (Yu et al., 2019) is an entity-level sentiment analysis method based on LSTM. 2) **TomBERT** (Yu and Jiang, 2019) applies BERT to obtain aspect-sensitive textual representations. 3) **CapTrBERT** (Khan and Fu, 2021) translates images into text and construct an auxiliary sentence for fusion. 4) **JML** (Ju et al., 2021) is the first joint model for MABSA with an auxiliary cross-modal relation detection module. 5) **VLP-MABSA** (Ling et al., 2022) performs five task-specific pretraining tasks to model aspects, opinions, and alignments. 6) **CMMT** (Yang et al., 2022) implements a gate to control the multimodal information contributions during inter-modal interactions. 7) **AoM** (Zhou et al., 2023) introduces an aspect-oriented network designed to reduce visual and textual distractions from complex image-text interactions. 8) **Emotion-LLaMA** (Cheng et al., 2024). 9) **AETS** (Zhu et al., 2025) improves multimodal sentiment analysis by enhancing aspects and simplifying text. 10) **Qwen3-VL-8B-Thinking** (Yang et al., 2025a)

## D Implementation Details

### D.1 Hyperparameters in Multi-Task Learning

In our multi-task learning setup, we assign weights of 0.8 and 0.2 to the CoT (Chain-of-Thought) generation task and the sentiment classification task, respectively. This design is motivated by the following considerations:

- **Task complexity:** CoT generation involves structured reasoning and belongs to a class of complex sequence generation tasks, which are more difficult to train and typically incur higher loss values. In contrast, sentiment

classification is a relatively simple three-way classification task. Therefore, assigning a higher weight to CoT generation encourages the model to focus more on learning reasoning capabilities.

- **Convergence and gradient sensitivity:** As shown in Figure 4, when the loss weights of the CoT generation task are set to 0.8 and 0.2, the model exhibits a significant difference in converged loss: specifically, the loss differs by approximately 4x on the MVSA-Single dataset (0.0134 vs. 0.0033) and by approximately 12x on the Twitter-2015 dataset (0.0084 vs. 0.0007). In contrast, for the sentiment classification task, the model’s converged loss remains largely consistent under different loss weight settings. Preliminary experiments show that the CoT task converges more slowly and is more sensitive to gradient fluctuations. Increasing its loss weight helps amplify gradient signals and improves training stability and task performance.
- **Empirical validation:** We experimented with different weight configurations (e.g., {0.5, 0.5}, {0.2, 0.8}) and observed that assigning lower weights to the CoT task led to slower loss reduction and decreased classification accuracy. In contrast, the {0.8, 0.2} setting consistently yielded better performance on both the validation and test sets.

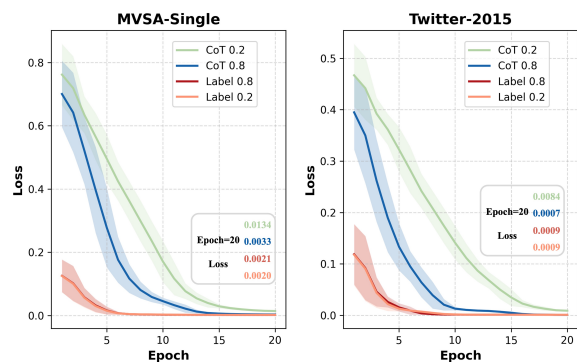


Figure 4: Effect of Loss Weight on Convergence for CoT Generation and Sentiment Classification Tasks.

This weighting scheme also reflects the task balancing principle proposed by CoTBal (Dai et al., 2024), which emphasizes that in multi-task scenarios, loss weights should be adaptively assigned based on task complexity and learning dynamics to

enhance main-task optimization and overall model performance.

## D.2 Hyperparameter in Knowledge Distillation

We set the hyperparameter  $\lambda$  to 0.3, following the empirical practices in prior work (Lee et al., 2024), which achieve a good balance between stable training and effective knowledge transfer from the teacher model.

## D.3 LoRA Configuration

In all experiments, we adopt Low-Rank Adaptation (LoRA) for parameter-efficient fine-tuning of the multimodal models, rather than updating all parameters. In our implementation, the LoRA rank is set to  $r = 16$  and the scaling factor is set to  $\alpha = 32$ . This strategy allows the models to adapt effectively to downstream tasks under limited computational resources, while significantly reducing the number of trainable parameters. The proportion of trainable parameters for each model is summarized in Table 10, and further implementation details are available in the anonymous code repository.

Model	Total Parameters	Trainable	Ratio (%)
Qwen3-VL-8B	8,810,770,672	43,646,976	0.4954
Qwen3-VL-2B	2,144,964,608	17,432,576	0.8127
Qwen2.5-VL-7B	8,339,756,032	47,589,376	0.5706
Qwen2.5-VL-3B	3,791,775,744	37,152,768	0.9798

Table 10: Total and trainable parameter counts under LoRA fine-tuning.

## E Details of the Adaptive Retry Controller (ARC)

Since large multimodal models are still prone to hallucinations in open-ended generation—such as producing content that deviates from the task specification or violating the required output format—the predicted sentiment label cannot always be reliably extracted from a single response. To mitigate this issue, we introduce an Adaptive Retry Controller (ARC) in the inference stage. Whenever the initially generated response does not conform to the expected format or fails to yield a valid sentiment label, ARC automatically triggers a retry process, prompting the model to regenerate the response for the same input sample. This process is repeated until a valid output is obtained or a predefined maximum number of retries is reached.

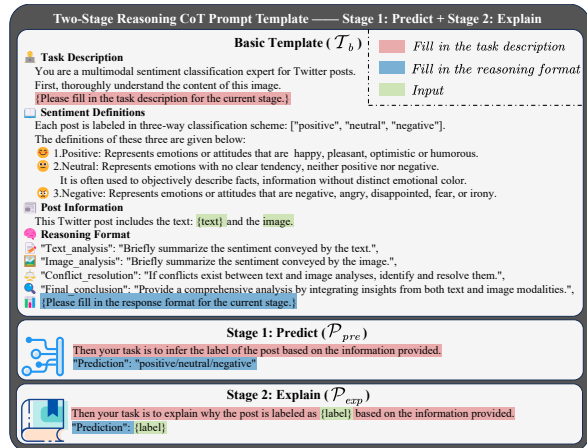


Figure 5: Two-stage reasoning prompt template.

In our implementation, the maximum number of retries is set to three. This choice strikes a balance between improving label extraction robustness and controlling computational overhead. Empirically, approximately 79% of the samples obtain a valid prediction within one or two generations, while only a small fraction require the full retry budget. As a result, ARC substantially reduces cases of missing or incorrectly parsed predictions, leading to a more stable and reliable inference pipeline without incurring prohibitive additional cost.

## F Two-stage reasoning prompt template

The prompt template for the two-stage multimodal Chain-of-Thought (CoT) is illustrated in Figure 5. By leveraging the basic template  $\mathcal{T}_b$ , the framework assigns a multimodal expert role to the model and prescribes a four-step structured reasoning process—comprising text analysis, image analysis, conflict resolution, and final conclusion—to explicitly address emotional correlations and conflicts between text and image modalities. Specifically, the execution is bifurcated into two stages: Prediction ( $\mathcal{P}_{pre}$ ) and Explanation ( $\mathcal{P}_{exp}$ ). In the first stage, the model is required to autonomously derive sentiment labels based on multimodal information to evaluate its inherent reasoning capabilities. In the second stage, the model is guided to backwardly construct logical justifications grounded in known labels, thereby generating deep-seated sentiment explanations. This two-stage decoupled design not only ensures the logical rigor of the reasoning chain but also substantially enhances the quality and interpretability of the generated data through the CoT mechanism.

## G Robustness of MulCoT-RD

For the Flan-T5 series. We utilize MiniCPM-o-2.6 (Team, 2025) to generate image captions, converting multimodal inputs to text-only format. Using the Flan-T5 architecture, we fine-tune both assistant and student models with full parameters, replicating the complete training pipeline including multimodal CoT enhancement, multi-task learning, and reasoning distillation. As shown in Figure 6 to 8, the Flan-T5-based models achieve strong performance despite having only 248M parameters, demonstrating the robustness and adaptability of MulCoT-RD across diverse backbone architectures.

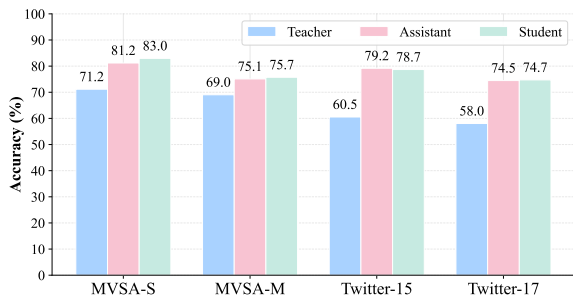


Figure 6: Accuracy comparison of teacher (GPT-3.5-Turbo), assistant (Flan-T5-Large with 783M parameters) and student (Flan-T5-Base) models.

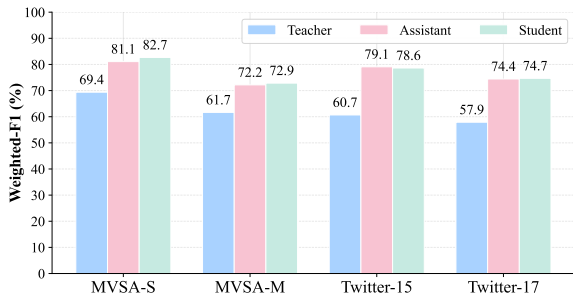


Figure 7: Weighted-F1 comparison of teacher (GPT-3.5-Turbo), assistant (Flan-T5-Large with 783M parameters) and student (Flan-T5-Base) models.

## H Efficiency of MulCoT-RD

From the comparison of the data in Tables 11 and Table 12, Qwen3-VL-2B demonstrates strong competitiveness under resource-constrained scenarios. On the MVSA-Single dataset, its GPU memory consumption is only 4.11 GB, approximately 13% of Emotion-LLaMA (31.61 GB), substantially lowering the hardware deployment barrier. In contrast, our student models, by optimizing TFLOPs, achieve significantly reduced end-to-end inference

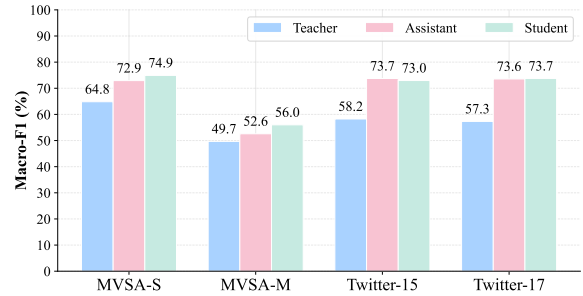


Figure 8: Macro-F1 comparison of teacher (GPT-3.5-Turbo), assistant (Flan-T5-Large with 783M parameters) and student (Flan-T5-Base with 248M parameters) models.

latency while maintaining minimal memory usage. These results clearly indicate that knowledge distillation using high-quality two-stage CoT data enables lightweight models to achieve a favorable balance of high accuracy, low latency, and compact footprint in multimodal sentiment analysis tasks, providing reliable support for real-time online monitoring and mobile deployment.

Model	Params (B)	TFLOPs -	Latency (ms)	Memory (GB)
Qwen2.5-VL-72B	73.41	38.24	17849.37	149.69
Emotion-LLaMA	14.58	4.11	11342.73	31.61
Qwen3-VL-8B	8.81	3.13	10708.24	16.61
Qwen3-VL-2B	2.14	3.15	8302.95	4.11
Qwen2.5-VL-7B	8.34	3.02	7002.99	15.81
Qwen2.5-VL-3B	3.79	1.32	8500.58	7.35

Table 11: Efficiency comparison on the MVSA-Single dataset.

Furthermore, we observe that Qwen2.5-VL-7B exhibits even lower inference latency than Qwen3-VL-2B and Qwen2.5-VL-3B, despite having a larger parameter count. This can be attributed to two factors: first, Qwen2.5-VL-7B has fewer LLM layers (28 vs. 36) (Bai et al., 2025); second, its larger hidden dimension increases the size of per-layer matrix operations, which allows modern GPUs to better saturate memory bandwidth and compute units, thereby achieving higher per-token computation efficiency (higher model bandwidth utilization) (Recasens et al., 2025).

## I Case Study

To validate MulCoT-RD's effectiveness, we present two illustrative cases in Figure 9. In case (a), D<sup>2</sup>R incorrectly predicts sentiment by overrelying on surface-level positive terms like "popular" and "bi-partisan" while missing the emotional shift from



Task	Coarse-Grained MSA	Fine-Grained MSA
Image		
Text	(a) RT @nationaljournal: Supreme Court term limits are popular, bipartisan, and hopeless.	(b) RT @tunkuv : Wedding photo becomes symbol of South Korea <b>MERS</b> scare
Previous SOTA	Method: D <sup>2</sup> R Prediction: Positive <b>X</b>	Method: AoM Prediction: Negative <b>X</b>
MulCoT-RD (ours)	<p><b>Sentiment Reasoning</b></p> <p>"Text_analysis": "The text discusses the popularity and bipartisan support of Supreme Court term limits but calls them 'hopeless' indicating negative sentiment about feasibility."</p> <p>"Image_analysis": "The image shows a large crowd outside the Supreme Court, suggesting public concern and urgency."</p> <p>"Conflict_resolution": "Despite some positive points, the overall tone is negative due to 'hopeless.' The crowd image reinforces public frustration or disappointment."</p> <p>"Final_conclusion": "The text's negative view aligns with the image's portrayal of public concern, jointly conveying an overall negative sentiment."</p>	<p><b>Sentiment Reasoning</b></p> <p>"Text_analysis": "The text says the wedding photo symbolizes the South Korea MERS scare but does not express positive or negative sentiment toward MERS, focusing on cultural significance with a neutral tone."</p> <p>"Image_analysis": "The image shows people, including a newlywed couple, wearing masks possibly related to MERS, but it conveys no clear sentiment and maintains a neutral tone."</p> <p>"Conflict_resolution": "There is no conflict; both text and image present neutral, factual observations without emotional expression."</p> <p>"Final_conclusion": "Overall, the text and image provide objective, neutral information about the symbolic meaning and health measures, supporting a neutral classification."</p>
	<p>"Prediction": "Negative" <b>✓</b></p> <p><b>Sentiment Classification</b></p>	<p>"Prediction": "Neutral" <b>✓</b></p> <p><b>Sentiment Classification</b></p>

Figure 9: Visualization of two samples, using the MulCoT-RD architecture with ID 4 from Table 2.

Model	Params (B)	TFLOPs -	Latency (ms)	Memory (GB)
Qwen2.5-VL-72B	73.41	38.41	17524.94	149.76
Emotion-LLaMA	14.58	4.31	11582.16	31.66
Qwen3-VL-8B	8.81	3.39	9930.38	16.62
Qwen3-VL-2B	2.14	3.60	8710.76	4.12
Qwen2.5-VL-7B	8.34	3.21	6209.65	15.83
Qwen2.5-VL-3B	3.79	1.44	8113.90	7.37

Table 12: Efficiency comparison on the Twitter-2015 dataset.

reasoning and sentiment classification, enabling comprehensive modeling of intra-modal and cross-modal sentiment reasoning.

1223  
1224  
1225

1213 the word "hopeless" which establishes a negative  
1214 tone. MulCoT-RD successfully captures this reversal.  
1215 In case (b), the AoM misclassifies sentiment  
1216 for the aspect term "MERS" by focusing on superficially  
1217 negative words like "scare", leading to misinterpretation.  
1218 MulCoT-RD effectively distinguishes between author stance  
1219 (factual reporting) and content sentiment, producing correct  
1220 predictions. This superior performance stems from our  
1221 multi-task learning mechanism that integrates CoT  
1222