

Emotion-Coherent Reasoning for Multimodal LLMs via Emotional Rationale Verifier

Hyeongseop Rha¹, Jeong Hun Yeo¹, Yeonju Kim¹, Yong Man Ro^{1*}

¹Integrated Vision and Language Lab, KAIST, South Korea
ryool.1832@kaist.ac.kr, sedne246@kaist.ac.kr, yeonju7.kim@kaist.ac.kr, ymro@kaist.ac.kr

Abstract

The recent advancement of Multimodal Large Language Models (MLLMs) is transforming human-computer interaction (HCI) from surface-level exchanges into more nuanced and emotionally intelligent communication. To realize this shift, emotion understanding becomes essential allowing systems to capture subtle cues underlying user intent. Furthermore, providing faithful explanations for predicted emotions is crucial to ensure interpretability and build user trust. However, current MLLM-based methods often generate emotion explanations that diverge from the target labels and sometimes even contradict their own predicted emotions. This inconsistency poses a critical risk for misunderstanding and erodes reliability in interactive settings. To address this, we propose a novel approach: the Emotional Rationale Verifier (ERV) and an Explanation Reward. Our method guides the model to produce reasoning that is explicitly consistent with the target emotion during multimodal emotion recognition without modifying the model architecture or requiring additional paired video–description annotations. Our method significantly improves faithful explanation–prediction consistency and explanation emotion accuracy on the MAFW and DFEW datasets. Through extensive experiments and human evaluations, we show that our approach not only enhances alignment between explanation and prediction but also empowers MLLMs to deliver emotionally coherent, trustworthy interactions, marking a key step toward truly human-like HCI systems.

1 Introduction

As multimodal large language models (MLLMs) (Chen et al. 2024; Li et al. 2024; Wang et al. 2024; Chu et al. 2024; Cheng et al. 2024b) continue to evolve, they are reshaping human-computer interaction (HCI) by enabling systems to interpret cross-modal inputs (face video and speech audio), perform contextual reasoning, and generate human-like responses (Fei et al. 2024; Park et al. 2024; Kim et al. 2025). To support truly natural and emotionally aware communication, such systems must go beyond surface-level responses and capture the nuance, intent, and emotional context of user inputs. Among these capabilities, emotion understanding plays a central role in facilitating deeper, more meaningful interaction—particularly in applications such as psychological counseling (Schmidgall et al. 2024; Lee et al. 2024), educa-

tional guidance (Schutz and Pekrun 2007), and empathetic dialogue systems (Fei et al. 2024; Zhang et al. 2025).

While conventional multimodal emotion recognition (MER) approaches have treated emotion understanding as a classification task (Zhang et al. 2024b), single-label predictions often fail to capture the complexity of human affect conveyed through face video and speech audio. To address this limitation, recent work has turned to *descriptive emotion understanding*, leveraging MLLMs to generate natural language explanations that justify predicted emotions (Lian et al. 2023; Cheng et al. 2024a; Lian et al.; Yang et al. 2025a). This shift is especially valuable in HCI scenarios, where users not only expect accurate emotional classification from visual and auditory cues but also seek coherent and human-like justifications to build trust and ensure transparency.

To enhance emotion reasoning, a recent study leverages Reinforcement Learning (RL) and Chain-of-Thought (CoT) prompting to train MLLMs (Zhao, Wei, and Bo 2025). Although this method improves classification accuracy, it fails to consistently produce emotionally coherent explanations. This limitation may be attributable to the nature of the datasets used in current training pipelines. While recent method leverages a variety of emotion datasets, such as EMER (Lian et al. 2023), MERR (Cheng et al. 2024a), and MER-Caption (Lian et al. 2025), only a small fraction of these include human-annotated rationales. The majority of training relies on larger datasets without explanation labels, where reward signals focus solely on the correctness of predicted emotions, regardless of whether the generated rationale is emotionally coherent. As a result, models may generate emotionally incoherent explanations without penalty, leading to a misalignment between predictions and rationales, as shown in Figure 1.

In this paper, we propose a novel framework for generating emotionally coherent explanations, even when trained on emotion datasets lacking human-annotated rationales. Our approach is based on the key assumption that LLMs possess sufficient linguistic and emotional knowledge to evaluate the consistency between a given emotion and a generated explanation. However, incorporating LLMs directly into the evaluation process can be computationally expensive and impractical in deployment settings. To address this, we distill the knowledge of the LLM into a lightweight model, referred to as the Emotional Rationale Verifier (ERV), which effi-

*Corresponding author.

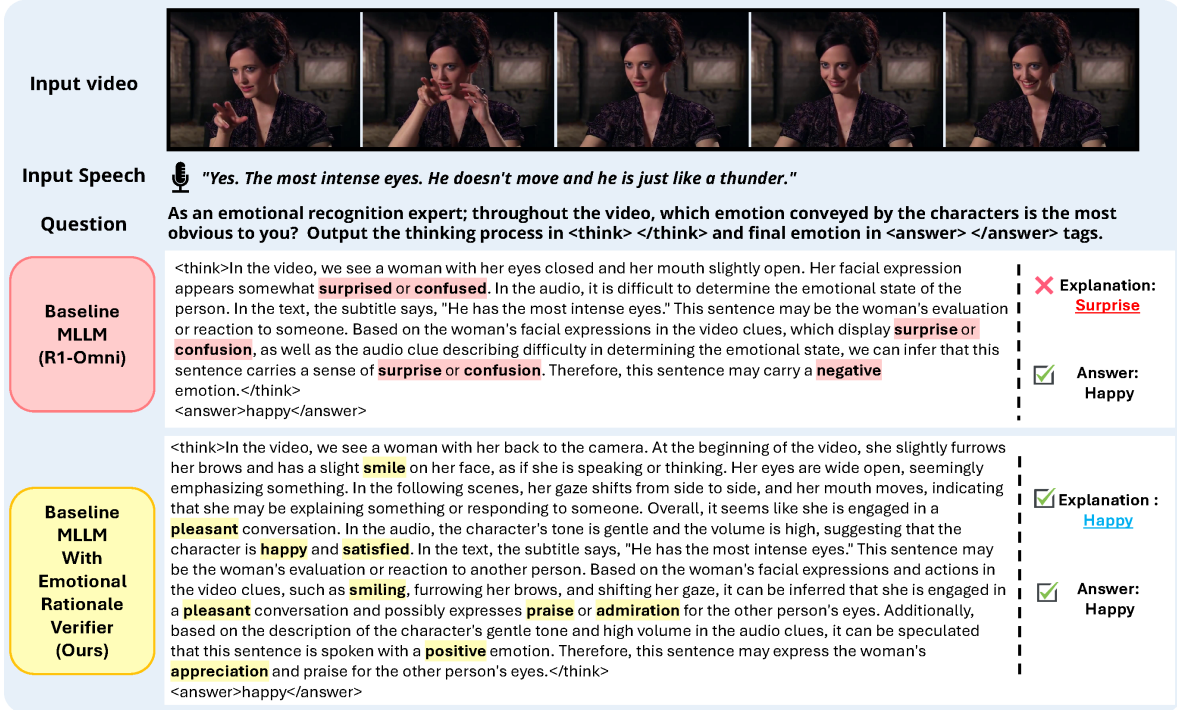


Figure 1: Illustrative comparison between the baseline and our model on explanation–emotion alignment. While both correctly predict the final emotion, the baseline fails to generate an emotionally coherent explanation, highlighting its misalignment.

ciently assesses the alignment between an explanation and the target emotion. Building on this verifier, we integrate ERV into a reward-driven training pipeline. In this setup, the model receives an explanation reward based on how well its generated explanation aligns with the GT emotion as judged by ERV. To further improve reward fidelity, we decompose each explanation into individual sentences, filter out emotionally neutral ones (e.g., those describing background or appearance), and assign rewards based on the proportion of emotionally salient sentences that correctly reflect the target emotion.

Furthermore, to evaluate the emotional coherence, we introduce a set of novel metrics: Explanation Emotion Accuracy (EEA), which measures alignment between the explanation and the target emotion, Explanation–Prediction Consistency (EPC), which quantifies the internal consistency between the explanation and predicted emotion, and Faithful Consistency Rate (FCR), a strict metric that requires all three components, explanation, predicted emotion, and target emotion, to agree. These metrics provide deeper insight into the emotional coherence of generated explanations beyond standard classification accuracy.

Experimental results demonstrate that our approach preserves emotion classification performance while improving the consistency between generated explanations and predicted emotions, as measured by the three proposed evaluation metrics. Moreover, human evaluation shows that the model produces more coherent and emotionally aligned reasoning, thereby supporting more reliable and interpretable

interactions in HCI applications.

Our key contributions are as follows:

- **Emotionally consistent explanation framework:** We propose a novel framework that improves the emotional coherence of explanations without modifying the model architecture or requiring labeled explanation data by leveraging CoT reasoning and reward-based learning.
- **Emotional Rationale Verifier (ERV):** We introduce ERV, a lightweight verifier distilled from an LLM, which assesses the alignment between generated explanations and target emotions, and serves as a reward function to guide explanation generation.
- **New metrics for emotional coherence evaluation:** We propose three novel metrics—EEA, EPC, and FCR—that enable fine-grained evaluation of emotional alignment and consistency in model-generated explanations.

2 Related Works

2.1 Multimodal Emotion Recognition (MER)

The task of Multimodal Emotion Recognition (MER) aims to understand human emotional states by analyzing diverse multimodal cues present in expressive scenes. From audio and visual cues to text, body, and physiological signals, numerous studies (Kossaifi et al. 2017; Zadeh et al. 2018; Poria et al. 2018; Busso et al. 2008; Perepelkina, Kazimirova, and Konstantinova 2018; Liu et al. 2022; Jiang et al. 2020; Livingstone and Russo 2018; Lee et al. 2019; Wang et al. 2022) have proposed datasets consisting of emotion-related modalities

and corresponding emotion labels. With those datasets, traditional approaches have primarily relied on modality-specific encoders (Zhang et al. 2017, 2021; Zhao and Liu 2021; Sun et al. 2023; Chumachenko, Iosifidis, and Gabbouj 2024) for feature extraction and fusion strategies (Narayan et al. 2024; Sun et al. 2024), achieving solid performance in emotion classification but falling short in terms of explainability.

However, unlike clear object classification, emotion recognition suffers from inherent ambiguity, as emotion labels alone often fail to fully capture an individual’s state. To address the challenge of unreliable labeling caused by the limited expressivity of fixed label mappings, multimodal emotion descriptive dataset EMER (Lian et al. 2023) has been proposed. This dataset not only requires models to predict emotion labels from visual, audio, and textual cues in videos, but also to generate descriptive explanations of the emotional states. Additionally, with the rise of Large Language Models (LLMs) (Touvron et al. 2023; Team 2024; Grattafiori et al. 2024), leveraging the natural language generation power, recent instruction-tuned multimodal LLMs approaches (Xing et al. 2024; Cheng et al. 2024a; Lian et al.; Yang et al. 2025a; Hu et al. 2025; Zhao et al. 2025a) move beyond classification toward emotion reasoning to enable interpretable, context-aware emotion understanding.

However, as these descriptions are based on human annotations, dataset sizes are limited, and scaling up emotional descriptive datasets is constrained by annotation costs. To address this issue, recent studies (Cheng et al. 2024a; Lian et al. 2024; Yang et al. 2025a) have investigated using the generative power of MLLMs to automatically augment emotional descriptive datasets. Building on these description-rich emotion datasets, many studies (Cheng et al. 2024a; Lian et al. 2024) have explored interpreting and generating explanations grounded in video and audio cues from samples. Furthermore, recent architectural advancements (Zhao et al. 2025b; Lian et al.) have contributed to improved recognition accuracy. However, these approaches typically rely on supervised fine-tuning (SFT), which inherently limits generalization due to its dependence on the scale and diversity of annotated emotional descriptions.

2.2 GRPO-Based RL with Verifiable Rewards

Recent work has shown that RL post-training can substantially sharpen the CoT abilities of MLLMs. In particular, Group Relative Policy Optimization (GRPO) (Shao et al. 2024) combined with Verifiable Reward (Guo et al. 2025) achieves strong generalization even with limited data and fixed model backbones, excelling in math and coding benchmarks.

Motivated by this success, GRPO has been adapted to a wide range of multimodal tasks, including emotion recognition (Zhao, Wei, and Bo 2025), general vision tasks (Chen et al. 2025a), object detection (Liu et al. 2025; Shen et al. 2025; Park et al. 2025), and the video domain (Li et al. 2025; Lee et al. 2025; Xing et al. 2025), often designing reward functions tailored to their specific task objectives. However, these approaches predominantly focus on optimizing the final answer, paying little attention to the quality and consistency of the intermediate reasoning steps, a degradation problem

highlighted by (Wei et al. 2025). To address this, recent studies (Chen et al. 2025b; Yang et al. 2025b) have introduced consistency-aware methodologies to improve reasoning and answer alignment. Yet, those consistency problems have not yet been explored in inherently ambiguous domains such as MER. We therefore introduce a novel reward applied with GRPO training that explicitly enforces coherence between the generated explanation and the target emotion.

3 Method

Our primary objective is to develop a Multimodal Large Language Model (MLLM) capable of emotion understanding that not only accurately predicts emotion labels from face video and speech audio but also generates consistent and faithful natural language explanations aligned with these predictions. Building upon recent advancements demonstrating the effectiveness of Reinforcement Learning (RL)-based training for improving both emotion recognition accuracy and reasoning in MLLM (Zhao, Wei, and Bo 2025), we integrate a novel Emotional Rationale Verifier (ERV) module into an established RL framework. This integration allows us to directly evaluate and enhance the emotional consistency of the generated explanations. Specifically, we extend the HumanOmni model architecture (Zhao et al. 2025b) and Group Relative Policy Optimization (GRPO) (Guo et al. 2025), a reward-driven RL training strategy showing strong performance in various multimodal tasks (Liu et al. 2025; Chen et al. 2025a; Zhao, Wei, and Bo 2025). While HumanOmni effectively unifies multimodal features and GRPO enhances Chain-of-Thought (CoT) reasoning, existing approaches often struggle with emotionally incongruent explanations despite accurate label predictions, as highlighted in our introduction.

To precisely address this critical misalignment, our proposed ERV module assesses whether a generated explanation coherently supports the target emotion. Based on this assessment, we define an additional reward signal, the Explanation Reward (R_E), which encourages the MLLM to produce more faithful and emotion-aligned rationales during RL fine-tuning. The overall training pipeline, detailing the ERV module and its integration into the GRPO framework, is visually represented in Figure 2.

3.1 Model Architecture and Training Overview

Our methodology leverages the HumanOmni architecture as its foundational MLLM. Given an input video V , audio A , and task prompt P , the model first processes the raw multimodal inputs using pretrained visual and audio encoders. These modality-specific features X_v (visual) and X_a (audio) are then projected into a shared embedding space and concatenated with the prompt embedding (X_p). The resulting comprehensive fused representation (X_m) is then fed into the Large Language Model (LLM), denoted as ψ , which autoregressively generates the output sequence (o):

$$X_m = \text{Concat}(X_v, X_a, X_p), \quad o = \psi(X_m) \quad (1)$$

The training of this model follows a robust two-stage strategy: an initial Supervised Fine-Tuning (SFT) (Zhang et al. 2024a) phase, followed by RL using GRPO. This progressive approach ensures the model first learns fundamental reason-

ing structures and then refines its explanation generation based on explicit reward signals.

Initializing Basic Emotion Reasoning via SFT The initial training phase involves SFT on a high-quality, albeit relatively small, dataset named EMER (Lian et al. 2023). This dataset consists of (V, A) pairs annotated with precise emotion labels and corresponding descriptive explanations. The primary objective of this phase is to enable the MLLM to grasp the basic structural relationships between multimodal inputs, emotion labels, and their textual explanations. It specifically helps the model learn to produce outputs in the required format $\langle \text{think} \rangle \dots \langle \text{think} \rangle \langle \text{answer} \rangle \dots \langle \text{answer} \rangle$. This SFT phase serves as a crucial initialization step for our policy model, denoted as π_θ (i.e., the trainable MLLM). During SFT, the model is optimized to maximize the likelihood of generating the target output sequence given the input. The training objective is defined by the negative log-likelihood loss:

$$\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{SFT}}} \left[\sum_{t=1}^{|y|} \log \pi_\theta(y_t \mid x, y_{<t}) \right] \quad (2)$$

Here, (x, y) represents the input-output pair sampled from the SFT dataset \mathcal{D}_{SFT} , where x represents the input triplet (V, A, P) and y is the corresponding explanation.

Enhancing Explanation Coherence through Explanation-Guided Reinforcement Learning Following the SFT phase, we transition to RL using GRPO to enhance the model’s reasoning capability and, critically, the consistency of its explanations. The SFT-trained model π_θ functions as the learnable policy, while a frozen duplicate, π_{ref} , acts as the reference model. For each input (V, A, P) , both the policy model and the reference model generate G candidate output sequences, $\{o_1, o_2, \dots, o_G\}$. Each generated sequence is then evaluated by a comprehensive reward function, which we detail in the next section. The policy π_θ is subsequently optimized to maximize the expected reward R , while simultaneously being regularized to remain close to the behavior of π_{ref} . The optimization objective of GRPO is defined as follows:

$$\max_{\pi_\theta} \mathbb{E}_{o \sim \pi_\theta(X)} [R(o) - \beta \cdot \text{KL}(\pi_\theta(o|X) \parallel \pi_{\text{ref}}(o|X))] \quad (3)$$

While the GRPO framework is highly effective for reward-driven optimization, it traditionally lacks an explicit mechanism to guarantee the consistency between the predicted emotion and its generated explanation. To overcome this inherent limitation and directly address the core challenge identified in our introduction, we introduce the Emotional Rationale Verifier (ERV) and its associated reward (R_E) as a novel, crucial signal within the GRPO training loop.

3.2 Emotional Rationale Verifier (ERV): Design and Training

We define an explanation-related reward signal based on how well the explanation in the output sequence o reflects the target emotion. To assess the emotional content of the generated explanation E , we employ an auxiliary module capable of classifying emotion labels from textual descriptions.

Although existing datasets such as EMER (Lian et al. 2023) and MERR-fine (Cheng et al. 2024a) provide pairs of emotional descriptions and emotion labels, their limited sizes (332 and 4,487 samples, respectively) are insufficient for our setting. This is particularly problematic because the generated descriptions often contain not only emotional expressions but also contextual details about people, environments, and temporally evolving emotional states. In addition, the emotion label distribution is imbalanced, as shown in Appendix B.1. As a result, existing datasets lack the diversity and coverage necessary to effectively train the Emotional Rationale Verifier (ERV).

To address the limitations of existing datasets in terms of both size and class imbalance, we construct a pseudo-labeled dataset consisting of emotional text descriptions paired with emotion labels. As shown in Figure 1, while R1-Omni (Zhao, Wei, and Bo 2025) does not always produce perfectly emotion-aligned outputs, its generated descriptions are often emotionally plausible. Leveraging this property, we generate textual explanations from the training sets of DFEW (Jiang et al. 2020) and MAFW (Liu et al. 2022), and use a closed-source LLM (GPT-4.1) to assign up to two representative emotion labels to each description. This process yields a pseudo-labeled dataset of 20K emotional text-label pairs, which is used to train our ERV module. To further mitigate the class imbalance, we augment underrepresented categories, namely *disgust*, *fear*, *contempt*, *disappointment*, *neutral*, and *helplessness*, by generating additional emotional descriptions using GPT-4.1 (see Figure 7). We provide few-shot exemplars for each of these emotion labels as prompts, and synthesize new descriptions such that each of these categories contains approximately 1K examples in the final dataset.

To train the ERV on the constructed dataset, we use the RoBERTa model trained on the GoEmotions dataset (Demszky et al. 2020)¹ as a backbone. We fine-tune this model to predict emotion labels from text descriptions using our 23K-sample dataset, which includes both human-annotated and GPT-generated pseudo-labeled examples. The model predicts the emotion label based on the output of the $[\text{CLS}]$ token, which encodes the global semantics of the input description. These predictions are then used to compute the Explanation Reward for each generated explanation E_i .

3.3 Explanation Reward

With the proposed ERV, we define an explanation reward R_E , which is assigned in proportion to how well the generated explanation E reflects the GT emotion e_{gt} . However, a long narrative often mixes various types of content such as scene background, appearance details, and affective cues, making it unreliable to evaluate the entire paragraph as a whole. Sentences that merely describe appearance or context can dilute the emotional signal and mislead ERV.

To address this, we design the explanation reward R_E with two key components. We illustrate the reward computation using the i th output o_i from the policy model π_θ , as shown in Figure 2(b). **1) Sentence-level and multi-label verifica-**

¹https://huggingface.co/SamLowe/roberta-base-go_emotions

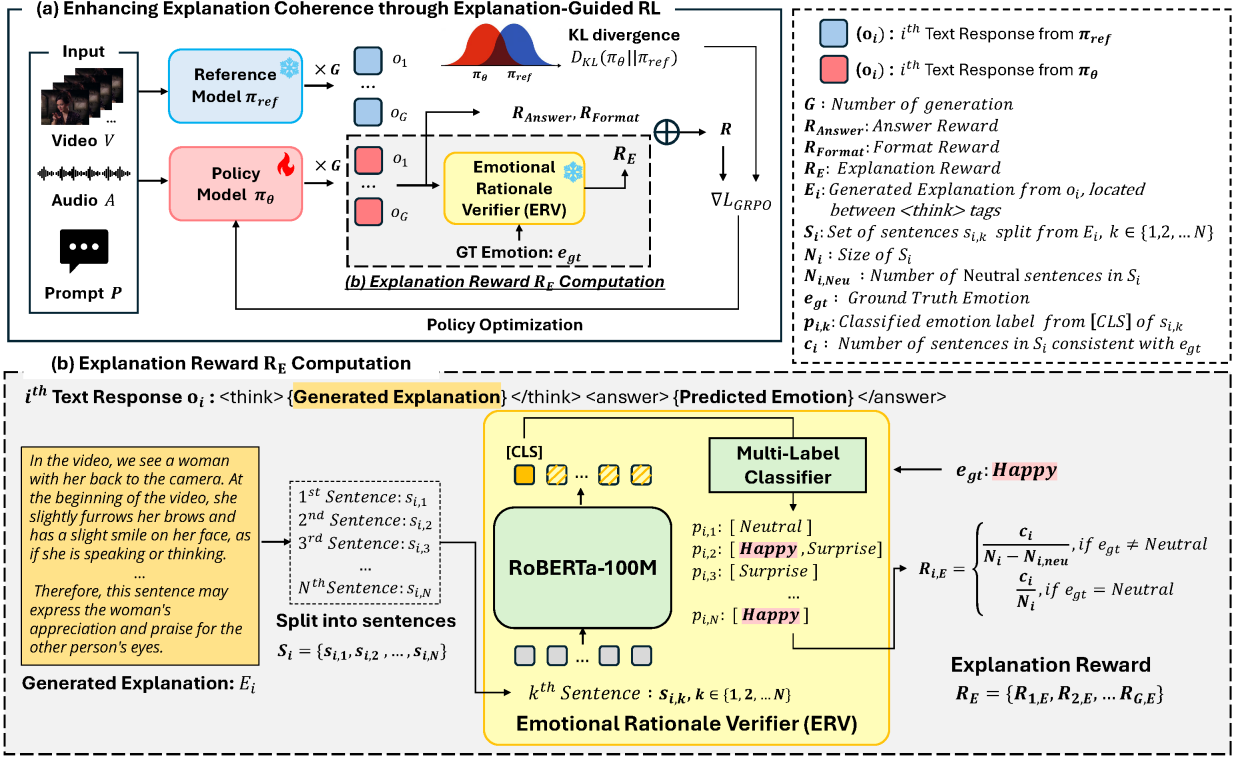


Figure 2: **(a) Enhancing explanation coherence using Explanation-Guided Reinforcement Learning (RL).** The policy model π_θ generates G responses o_1, \dots, o_G , and the Emotional Rationale Verifier (ERV) assigns an explanation reward R_E to each response. **(b) Explanation Reward R_E Computation.** For each response o_i , the corresponding explanation E_i is extracted and evaluated to produce its reward $R_{i,E}$. Collectively, $R_E = \{R_{1,E}, R_{2,E}, \dots, R_{G,E}\}$ denotes the set of explanation rewards for the G generated responses.

tion. For each generated sentence $s_{i,k} \in E_i$, ERV predicts a (possibly multi-label) emotion set $p_{i,k}$. This sentence-level prediction prevents non-emotive sentences from dominating the overall judgment. Furthermore, since a single sentence may convey multiple emotions, we apply a threshold τ to extract up to two emotion labels per sentence from the classifier. **2) Neutral-sentence filtering.** If ERV classifies $s_{i,k}$ as NEUTRAL, the sentence is excluded from the computation of $R_{i,E}$, based on the intuition that surface-level descriptions such as listing physical traits or environmental details provide little to no emotional signal and thus do not meaningfully contribute to emotion understanding.

Given the set of split sentences S_i , each sentence is evaluated by ERV. We count the number of sentences whose predicted emotions $p_{i,k}$ match the GT emotion e_{gt} as c_i , and denote the number of sentences classified as NEUTRAL as $N_{i,neu}$. The explanation reward $R_{i,E}$ is then computed as:

$$R_{i,E} = \begin{cases} \frac{c_i}{N_i - N_{i,neu}}, & \text{if } e_{gt} \neq \text{NEUTRAL}, \\ \frac{c_i}{N_i}, & \text{if } e_{gt} = \text{NEUTRAL}. \end{cases} \quad (4)$$

Final reward. In addition to the explanation reward R_E , we incorporate two auxiliary rewards *format reward* and *answer reward* to compute the final total reward R . The format reward R_{format} , following prior work such as RLVR, assigns

1 if the model output adheres to the required format (i.e., includes both $\langle \text{think} \rangle$ and $\langle \text{answer} \rangle$ tags), and 0 otherwise. The answer reward R_{answer} assigns 1 if the predicted emotion within the $\langle \text{answer} \rangle$ section matches the GT emotion label e_{gt} .

The final reward used in training is computed as:

$$R = R_E + R_{format} + R_{answer}. \quad (5)$$

By incorporating the explanation reward R_E obtained from ERV as an additional learning signal during GRPO training, we aim to enhance the consistency between the model’s reasoning process and the GT emotion.

4 Experimental Setup

4.1 Dataset

EMER (Lian et al. 2023) is a human-annotated emotional video dataset with 332 samples labeled with one of five emotions (*angry*, *sad*, *surprise*, *worried*, *happy*) and corresponding descriptions. We use EMER to train the HumanOmni model to generate explanations in a reasoning format.

MERR (Cheng et al. 2024a) provides 4,487 fine-grained pseudo-labeled video-caption pairs, generated via facial and audio analysis with LLM-guided annotation. We use this data to supervise our ERV, leveraging its diverse and component-level emotion cues.

Model	MAFW					DFEW				
	EEA (%)	FCR (%)	EPC (%)	WAR (%)	UAR (%)	EEA (%)	FCR (%)	EPC (%)	WAR (%)	UAR (%)
SFT-based Emotion MLLMs										
Human-Omni (Zhao et al. 2025b)	–	–	–	20.18	13.52	–	–	–	22.64	19.44
EMER-SFT	33.06	29.67	70.66	38.95	25.21	31.47	24.28	56.49	38.93	30.86
EMER-SFT (7B)	44.26	38.58	67.38	50.22	33.48	40.68	27.58	42.27	48.35	50.13
<i>Emotion-LLaMA (7B)</i>	–	–	–	–	–	31.70	21.39	42.14	48.71	38.54
<i>Emotion-LLaMA (7B) Fine-Tuned</i>	–	–	–	–	–	–	–	–	77.06	64.21
RL-based Emotion MLLMs										
R1-Omni (Zhao, Wei, and Bo 2025)	37.87	31.80	48.09	57.68	40.04	33.49	28.09	41.58	65.44	56.94
Baseline	32.41	27.27	40.98	56.83	38.66	32.37	27.57	40.16	66.17	56.10
<i>Baseline (7B)</i>	44.97	40.22	55.19	<u>65.36</u>	<u>47.27</u>	43.34	38.54	46.94	<u>77.09</u>	<u>73.73</u>
Ours										
Ours	48.93	43.57	68.58	57.49	39.12	45.82	40.21	58.63	67.28	56.25
Ours (7B)	<u>54.70</u>	<u>50.98</u>	<u>73.06</u>	65.19	47.01	<u>50.88</u>	<u>45.95</u>	<u>57.09</u>	75.85	72.15

Table 1: Evaluation of models on proposed metrics: Explanation Emotion Accuracy (EEA), Faithful Consistency Rate (FCR), Explanation–Prediction Consistency (EPC), and classification accuracy (WAR, UAR) on MAFW and DFEW datasets. Underline indicates the best score based on the 7B-sized model, while bold indicates the best score based on the 0.5B-sized model.

DFEW (Jiang et al. 2020) is a large-scale facial expression video dataset with 9,362 training and 2,342 test samples, labeled with 7 basic emotions: *angry, disgust, surprise, happy, sad, neutral, and fear* (official 5th split). We use DFEW in two ways: (1) generating textual explanations from training videos and assigning pseudo-labels via GPT-4.1 to augment ERV training, and (2) applying our reward-based RL scheme using EMER-pretrained models to produce `<think>` and `<answer>` outputs.

MAFW (Liu et al. 2022) is a multimodal emotion dataset with 7 basic and 4 compound categories (*contempt, helplessness, anxiety, disappointment*), comprising 7,341 training and 1,831 test samples (*set 5 split*). We apply the same pipeline as DFEW, pseudo-labeled explanation augmentation for ERV and RL training using our proposed rewards.

4.2 Evaluation Metrics

Emotion Recognition Accuracy Following prior works (Zhao, Wei, and Bo 2025), we employ Unweighted Average Recall (UAR) and Weighted Average Recall (WAR) to measure the model’s emotion classification performance.

Evaluating Emotional Coherence We evaluate the emotional coherence of generated explanations using three metrics: Explanation Emotion Accuracy (EEA), Explanation–Prediction Consistency (EPC), and Faithful Consistency Rate (FCR). These metrics capture different aspects of alignment among the explanation, predicted emotion, and ground-truth label. GPT-4.1-mini is used to recognize emotion in explanation. Figure 3 shows the prompt used in our evaluation. The detailed formulations of these metrics are provided in Appendix ??.

4.3 Implementation Details

We employ HumanOmni 0.5B model (Zhao et al. 2025b) as a backbone model. Our ERV is based on the RoBERTA-100M trained on the GoEmotions Dataset. Our training pipeline consists of three components: SFT, ERV module, and GRPO training. The main model is first fine-tuned on the EMER dataset for 5 epochs with a learning rate of $2e^{-5}$, a cosine

```

Read the reasoning content and respond
with the appropriate emotion in {Emotion
List}.
Reply with only the emotion word.
Reason: {Explanation} Answer Emotion:

```

Figure 3: Prompt used to evaluate the emotion conveyed in the generated explanation. {Emotion List} is a permutation of the GT emotion label according to the evaluation dataset. {Explanation} is E_i from output o_i .

scheduler, a warmup ratio of 0.03, and a batch size of 32. The ERV module is trained separately for 10 epochs with a batch size of 64 and a learning rate of $3e^{-5}$. Those are trained on 8 NVIDIA RTX 3090 GPUs. For GRPO training, we use a batch size of 16 with a gradient accumulation step of 2, generating $G = 16$ samples per input, while applying a KL divergence penalty with coefficient $\beta = 0.04$ and fixing the sampling temperature at 0.3 during inference. This stage is conducted on 8 NVIDIA A100 GPUs. Appendix E shows more training details.

5 Experimental Results

5.1 Main Results

To verify the effectiveness of our proposed Explanation reward in the quality of generated explanations, we evaluated on three proposed metrics: EEA, EPC, and FCR, with standard recognition accuracy WAR and UAR. Our model is compared with three models trained in different ways using the same multimodal emotion description dataset and HumanOmni architecture. (1) **EMER-SFT**: A supervised fine-tuned model trained on the EMER dataset for emotion reasoning, available in both 0.5B and 7B LLM versions. (2) **R1-Omni**: A model trained with the GRPO using answer and format rewards. (3) **Baseline**: Since R1-Omni lacks training details and relies on an in-house dataset, we re-trained EMER-SFT with GRPO on the answer and format rewards

using the same training dataset and hyperparameter settings as our method for a fair comparison.

As shown in Table 1, both EMER-SFT and R1-Omni exhibit low FCR, achieving only 29.67% and 31.80%, respectively. This indicates that although R1-Omni significantly improves emotion recognition accuracy on WAR, from 38.95% to 57.68% compared to EMER-SFT, 44.87% of its predictions (almost half) are not supported by emotionally aligned explanations. Furthermore, its ability to independently produce explanations aligned with the target emotions improves only marginally from 33.06% to 37.87%. In contrast, our model trained with an Explanation reward not only achieves comparable recognition performance to R1-Omni but also demonstrates substantial improvements in both FCR and EEA. FCR increases from 29.67% to 43.57% and also EEA increases from 33.06% to 48.93%. A similar tendency is observed on the DFEW, where FCR improves from 24.28% to 40.21% and EEA improves from 31.47% to 45.82%.

We also compared our model with EMER-SFT (7B) and Emotion-LLaMA (7B), both of which use 7B-sized backbones but have different architectures. In the evaluation on the DFEW dataset, both models exhibited low performance, achieving EEA scores of 31.47% and 31.70%, and FCR scores of 24.28% and 21.39%, respectively. These results indicate that SFT training on existing multimodal video-reasoning pair datasets yields unreliable emotion reasoning performance regardless of the model architecture. In contrast, despite using a smaller 0.5B LLM size, our model outperforms both baselines in terms of FCR and EEA, while also achieving noticeably better emotion recognition accuracy. When applied to 7B-scale models, our method shows comparable recognition accuracy to Emotion-LLaMA fine-tuned solely for emotion label prediction, while achieving 45.95% FCR and 50.88% EEA.

About EPC, EMER-SFT shows significant performance compared to RL-based trained models. Namely, the final predicted answer is derived based on the generated explanations. However, when trained with GRPO on answer and format reward, this consistency deteriorates, as the model tends to focus solely on producing the final answer, leading to a degradation in EPC. In contrast, our proposed Explanation Reward enables the model to maintain this consistency, achieving comparable EPC performance to EMER-SFT (70.66% vs. 68.58% in MAFW and 56.49% vs. 58.63% in DFEW). This indicates that our Explanation Reward effectively guides the model to generate explanations and predictions in a coherent and aligned manner.

5.2 Ablation Study

Effectiveness of Emotional Rationale Verifier To verify our proposed training strategy on the ERV module, we removed the data augmentation stage utilizing GPT-4.1. (1) Baseline, without the ERV module, can’t reward on the generated explanation. It only focuses on the format and final answer. (2) Without a data augmentation strategy, training the ERV solely on the EMER and MERR-fine datasets boosts EEA from 32.41% to 46.88% and FCR from 27.27% to 42.62%, while keeping WAR and UAR unchanged. (3) Adding the augmented data to ERV further raises EEA +

Setting	EEA (%)	FCR (%)	EPC (%)	WAR (%)	UAR (%)
Baseline	32.41	27.27	40.98	56.83	38.66
Emotional Rationale Verifier (ERV)					
W/o Data Augmentation	46.88	42.62	65.46	56.78	38.60
Explanation Reward (R_E)					
W/o Neutral Filtering	46.28	42.13	64.70	57.16	38.67
W/o Sentence-Level Verification	47.54	43.17	68.36	56.78	38.66
Ours	48.93	43.57	68.58	57.49	39.12

Table 2: Ablation study on the Emotional Rationale Verifier (ERV) and Explanation Reward on MAFW

1.26 pp, FCR + 1.31 pp, EPC + 2.90 pp, and overall recognition accuracy. These results suggest that the ERV’s judgment capability is constrained by the size of its training data. When additional multimodal-emotion descriptions are unavailable, text-level data augmentation noticeably enhances the model’s multimodal reasoning ability.

Effectiveness of Explanation Reward To validate the effectiveness of our proposed Explanation Reward, we compare it against a variant trained without neutral filtering in the reward computation. That is, the reward does not take into account the proportion of non-emotional sentences (as described in Equation 4 for the neutral case). Without neutral filtering, the overall performance drops by 2.65 pp on EEA, 1.44 pp on FCR, 3.88 pp on EPC, 0.33 pp on WAR, and 0.45 pp on UAR. Furthermore, when sentence-level explanation verification is removed, we also observed performance degradation, with EEA and FCR dropping by 1.39 pp and 0.40 pp, respectively, and recognition metrics (WAR and UAR) also decreasing by 0.71 pp and 0.46 pp. As a result, when providing emotional guidance for explanation, it is effective to provide rewards according to emotional sentence proportion at the sentence level.

5.3 Human Evaluation about Explanation Quality

Through the main results and ablation studies, we confirmed that explanations from our model better align with both ground-truth and predicted emotions. However, their quality in terms of how well they explain the emotions in the videos had not been directly evaluated. To address this, we conducted a human study comparing our model’s explanations to those from R1-Omni. Specifically, we randomly selected 14 samples from each model (total 28), where the generated emotion matched the ground-truth. Twenty human raters evaluated each explanation on a 1–5 scale based on how well it explained the corresponding video’s emotion. For each sample, we recorded which model received the higher score and computed the average score across all samples. As shown in Table 3, our model achieved a 53.6% win rate and a higher average score (3.52 vs. 3.42), indicating that it produces more informative and emotionally grounded explanations than R1-Omni. (Survey details are provided in Appendix A.)

Evaluation	Ours	R1-Omni
Win Rate	53.6%	46.4%
Assessment Score	3.52	3.42

Table 3: Human evaluation comparison between our model and R1-Omni

Closed-source LLM	Model	EEA	FCR	EPC
GPT-4.1-mini	Baseline	32.41	27.27	40.98
	Ours	48.93	43.57	68.58
Gemini-Flash-2.5	Baseline	32.24	27.76	41.80
	Ours	48.52	42.73	65.96

Table 4: Comparison of Baseline and Ours across different closed-source LLMs in terms of EEA, FCR, and EPC.

5.4 Robustness of the Proposed Reasoning Evaluation Metrics

Since the evaluation of EEA, FCR, and EPC was conducted using a single closed-source LLM (GPT-4.1-mini), we further validated the robustness of the proposed metrics by performing additional evaluations with a different model family, *Gemini-2.5-Flash* (Comanici et al. 2025). As shown in Table 4, the evaluation on the MAFW for both the Baseline and our model exhibits the same tendency of improved reasoning quality. This demonstrates that our metrics are independent of the specific LLM used for evaluation and consistently capture the alignment between reasoning and the expressed emotion across different model families. The small absolute differences (0.17–2.62 percentage points, approximately 0.5–3.8% relative deviation) further confirm the robustness and stability of the evaluation across measurement settings. Moreover, in the human survey, participants were asked to judge the emotion conveyed by the model-generated reasoning, and their judgments showed a 92.9% agreement with those of GPT-4.1-mini (see Appendix A).

6 Conclusion

We proposed training an MLLM for the MER task utilizing our proposed ERV and Explanation Reward, aiming to enhance the coherence between emotion recognition and explanation. Compared to prior approaches, where coherence between predictions and explanations among correctly predicted samples remained around half (55.13% on MAFW and 42.92% on DFEW), our method significantly improves this alignment, achieving 75.79% and 59.77%. Furthermore, human evaluations confirmed that our approach produces more emotionally and contextually grounded emotion explanations for video data. As a result, our work enables emotionally coherent and trustworthy explanations alongside accurate emotion recognition, representing a significant advancement toward emotionally intelligent human–computer interaction.

References

- Busso, C.; Bulut, M.; Lee, C.-C.; Kazemzadeh, A.; Mower, E.; Kim, S.; Chang, J. N.; Lee, S.; and Narayanan, S. S. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42: 335–359.
- Chen, L.; Li, L.; Zhao, H.; Song, Y.; and Vinci. 2025a. R1-V: Reinforcing Super Generalization Ability in Vision-Language Models with Less Than \$3. <https://github.com/Deep-Agent/R1-V>. Accessed: 2025-02-02.
- Chen, Y.; Ge, Y.; Wang, R.; Ge, Y.; Cheng, J.; Shan, Y.; and Liu, X. 2025b. GRPO-CARE: Consistency-Aware Reinforcement Learning for Multimodal Reasoning. *arXiv preprint arXiv:2506.16141*.
- Chen, Z.; Wu, J.; Wang, W.; Su, W.; Chen, G.; Xing, S.; Zhong, M.; Zhang, Q.; Zhu, X.; Lu, L.; et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 24185–24198.
- Cheng, Z.; Cheng, Z.-Q.; He, J.-Y.; Wang, K.; Lin, Y.; Lian, Z.; Peng, X.; and Hauptmann, A. 2024a. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning. *Advances in Neural Information Processing Systems*, 37: 110805–110853.
- Cheng, Z.; Leng, S.; Zhang, H.; Xin, Y.; Li, X.; Chen, G.; Zhu, Y.; Zhang, W.; Luo, Z.; Zhao, D.; et al. 2024b. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*.
- Chu, Y.; Xu, J.; Yang, Q.; Wei, H.; Wei, X.; Guo, Z.; Leng, Y.; Lv, Y.; He, J.; Lin, J.; et al. 2024. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*.
- Chumachenko, K.; Iosifidis, A.; and Gabbouj, M. 2024. MMA-DFER: MultiModal Adaptation of unimodal models for Dynamic Facial Expression Recognition in-the-wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4673–4682.
- Comanici, G.; Bieber, E.; Schaekermann, M.; Pasupat, I.; Sachdeva, N.; Dhillon, I.; Blistein, M.; Ram, O.; and et al. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. *arXiv:2507.06261*.
- Demszky, D.; Movshovitz-Attias, D.; Ko, J.; Cowen, A.; Nemade, G.; and Ravi, S. 2020. GoEmotions: A Dataset of Fine-Grained Emotions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4040–4054.
- Fei, H.; Zhang, H.; Wang, B.; Liao, L.; Liu, Q.; and Cambria, E. 2024. EmpathyEar: An Open-source Avatar Multimodal Empathetic Chatbot. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, 61–71.
- Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

- Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Hu, Z.; Yuan, K.; Liu, X.; Yu, Z.; Zong, Y.; Shi, J.; Yue, H.; and Yang, J. 2025. Feallm: Advancing facial emotion analysis in multimodal large language models with emotional synergy and reasoning. *arXiv preprint arXiv:2505.13419*.
- Jiang, X.; Zong, Y.; Zheng, W.; Tang, C.; Xia, W.; Lu, C.; and Liu, J. 2020. Dfew: A large-scale database for recognizing dynamic facial expressions in the wild. In *Proceedings of the 28th ACM international conference on multimedia*, 2881–2889.
- Kim, Y.; Chung, J.; Kim, J.; Lee, S.; Lee, S.; Kim, J.; Yang, C.; and Yu, Y. 2025. Speaking Beyond Language: A Large-Scale Multimodal Dataset for Learning Nonverbal Cues from Video-Grounded Dialogues. *arXiv preprint arXiv:2506.00958*.
- Kossaifi, J.; Tzimiropoulos, G.; Todorovic, S.; and Pantic, M. 2017. AFEW-VA database for valence and arousal estimation in-the-wild. *Image and Vision Computing*, 65: 23–36.
- Lee, H.; Kim, J.; Kim, H.; and Ro, Y. M. 2025. ReFoCUS: Reinforcement-guided Frame Optimization for Contextual Understanding. *arXiv preprint arXiv:2506.01274*.
- Lee, J.; Kim, S.; Kim, S.; Park, J.; and Sohn, K. 2019. Context-aware emotion recognition networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, 10143–10152.
- Lee, S.; Mac Kim, S.; Kim, M.; Kang, D.; Yang, D.; Kim, H.; Kang, M.; Jung, D.; Kim, M.; Lee, S.; et al. 2024. Cactus: Towards Psychological Counseling Conversations using Cognitive Behavioral Theory. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 14245–14274.
- Li, F.; Zhang, R.; Zhang, H.; Zhang, Y.; Li, B.; Li, W.; Ma, Z.; and Li, C. 2024. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*.
- Li, X.; Yan, Z.; Meng, D.; Dong, L.; Zeng, X.; He, Y.; Wang, Y.; Qiao, Y.; Wang, Y.; and Wang, L. 2025. Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*.
- Lian, Z.; Chen, H.; Chen, L.; Sun, H.; Sun, L.; Ren, Y.; Cheng, Z.; Liu, B.; Liu, R.; Peng, X.; et al. 2025. AffectGPT: A New Dataset, Model, and Benchmark for Emotion Understanding with Multimodal Large Language Models. In *Forty-second International Conference on Machine Learning*.
- Lian, Z.; Chen, H.; Chen, L.; Sun, H.; Sun, L.; Ren, Y.; Cheng, Z.; Liu, B.; Liu, R.; Peng, X.; et al. 2025. Affectgpt: A new dataset, model, and benchmark for emotion understanding with multimodal large language models. *arXiv preprint arXiv:2501.16566*.
- Lian, Z.; Sun, H.; Sun, L.; Chen, L.; Chen, H.; Gu, H.; Wen, Z.; Chen, S.; Siyuan, Z.; Yao, H.; et al. 2024. Open-vocabulary multimodal emotion recognition: Dataset, metric, and benchmark.
- Lian, Z.; Sun, L.; Xu, M.; Sun, H.; Xu, K.; Wen, Z.; Chen, S.; Liu, B.; and Tao, J. 2023. Explainable multimodal emotion reasoning. *CoRR*.
- Liu, Y.; Dai, W.; Feng, C.; Wang, W.; Yin, G.; Zeng, J.; and Shan, S. 2022. Mafw: A large-scale, multi-modal, compound affective database for dynamic facial expression recognition in the wild. In *Proceedings of the 30th ACM international conference on multimedia*, 24–32.
- Liu, Z.; Sun, Z.; Zang, Y.; Dong, X.; Cao, Y.; Duan, H.; Lin, D.; and Wang, J. 2025. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*.
- Livingstone, S. R.; and Russo, F. A. 2018. The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PloS one*, 13(5): e0196391.
- Narayan, K.; VS, V.; Chellappa, R.; and Patel, V. M. 2024. Facexformer: A unified transformer for facial analysis. *arXiv preprint arXiv:2403.12960*.
- Park, S.; Kim, C.; Rha, H.; Kim, M.; Hong, J.; Yeo, J.; and Ro, Y. 2024. Let’s Go Real Talk: Spoken Dialogue Model for Face-to-Face Conversation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 16334–16348.
- Park, S.; Kim, H.; Kim, J.; Kim, S.; and Ro, Y. M. 2025. DIP-R1: Deep Inspection and Perception with RL Looking Through and Understanding Complex Scenes. *arXiv preprint arXiv:2505.23179*.
- Perepelkina, O.; Kazimirova, E.; and Konstantinova, M. 2018. RAMAS: Russian multimodal corpus of dyadic interaction for affective computing. In *International Conference on Speech and Computer*, 501–510. Springer.
- Poria, S.; Hazarika, D.; Majumder, N.; Naik, G.; Cambria, E.; and Mihalcea, R. 2018. Meld: A multimodal multi-party dataset for emotion recognition in conversations. *arXiv preprint arXiv:1810.02508*.
- Schmidgall, S.; Ziaei, R.; Harris, C.; Reis, E. P.; Jopling, J.; and Moor, M. 2024. AgentClinic: a multimodal agent benchmark to evaluate AI in simulated clinical environments. *CoRR*.
- Schutz, P. A.; and Pekrun, R. E. 2007. *Emotion in education*. Elsevier Academic Press.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Bi, X.; Zhang, H.; Zhang, M.; Li, Y.; Wu, Y.; et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Shen, H.; Liu, P.; Li, J.; Fang, C.; Ma, Y.; Liao, J.; Shen, Q.; Zhang, Z.; Zhao, K.; Zhang, Q.; et al. 2025. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*.
- Sun, L.; Lian, Z.; Liu, B.; and Tao, J. 2023. Mae-dfer: Efficient masked autoencoder for self-supervised dynamic facial expression recognition. In *Proceedings of the 31st ACM International Conference on Multimedia*, 6110–6121.
- Sun, L.; Lian, Z.; Liu, B.; and Tao, J. 2024. Hicmae: Hierarchical contrastive masked autoencoder for self-supervised

- audio-visual emotion recognition. *Information Fusion*, 108: 102382.
- Team, Q. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Wang, P.; Bai, S.; Tan, S.; Wang, S.; Fan, Z.; Bai, J.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; et al. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. *CoRR*.
- Wang, Y.; Sun, Y.; Huang, Y.; Liu, Z.; Gao, S.; Zhang, W.; Ge, W.; and Zhang, W. 2022. Ferv39k: A large-scale multi-scene dataset for facial expression recognition in videos. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 20922–20931.
- Wei, T.; Yang, Y.; Xing, J.; Shi, Y.; Lu, Z.; and Ye, D. 2025. GTR: Guided Thought Reinforcement Prevents Thought Collapse in RL-based VLM Agent Training. *arXiv preprint arXiv:2503.08525*.
- Xing, B.; Yu, Z.; Liu, X.; Yuan, K.; Ye, Q.; Xie, W.; Yue, H.; Yang, J.; and Kälviäinen, H. 2024. Emo-llama: Enhancing facial emotion understanding with instruction tuning. *arXiv preprint arXiv:2408.11424*.
- Xing, Z.; Hu, X.; Fu, C.-W.; Wang, W.; Dai, J.; and Heng, P.-A. 2025. Echoink-r1: Exploring audio-visual reasoning in multimodal llms via reinforcement learning. *arXiv preprint arXiv:2505.04623*.
- Yang, Q.; Bai, D.; Peng, Y.-X.; and Wei, X. 2025a. Omni-emotion: Extending video mllm with detailed face and audio modeling for multimodal emotion analysis. *arXiv preprint arXiv:2501.09502*.
- Yang, Q.; Yao, S.; Chen, W.; Fu, S.; Bai, D.; Zhao, J.; Sun, B.; Yin, B.; Wei, X.; and Zhou, J. 2025b. HumanOmniV2: From Understanding to Omni-Modal Reasoning with Context. *arXiv preprint arXiv:2506.21277*.
- Zadeh, A. B.; Liang, P. P.; Poria, S.; Cambria, E.; and Morency, L.-P. 2018. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2236–2246.
- Zhang, H.; Meng, Z.; Luo, M.; Han, H.; Liao, L.; Cambria, E.; and Fei, H. 2025. Towards multimodal empathetic response generation: A rich text-speech-vision avatar-based benchmark. In *Proceedings of the ACM on Web Conference 2025*, 2872–2881.
- Zhang, S.; Sun, P.; Chen, S.; Xiao, M.; Shao, W.; Zhang, W.; Liu, Y.; Chen, K.; and Luo, P. 2024a. Gpt4roi: Instruction tuning large language model on region-of-interest. In *European conference on computer vision*, 52–70. Springer.
- Zhang, S.; Tao, X.; Chuang, Y.; and Zhao, X. 2021. Learning deep multimodal affective features for spontaneous speech emotion recognition. *Speech Communication*, 127: 73–81.
- Zhang, S.; Yang, Y.; Chen, C.; Zhang, X.; Leng, Q.; and Zhao, X. 2024b. Deep learning-based multimodal emotion recognition from audio, visual, and text modalities: A systematic review of recent advancements and future prospects. *Expert Systems with Applications*, 237: 121692.
- Zhang, S.; Zhang, S.; Huang, T.; and Gao, W. 2017. Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching. *IEEE transactions on multimedia*, 20(6): 1576–1590.
- Zhao, F.; Li, M.; Xu, L.; Jiang, W.; Gao, J.; and Yan, D. 2025a. FaVChat: Unlocking Fine-Grained Facial Video Understanding with Multimodal Large Language Models. *arXiv preprint arXiv:2503.09158*.
- Zhao, J.; Wei, X.; and Bo, L. 2025. R1-omni: Explainable omni-multimodal emotion recognition with reinforcement learning. *arXiv preprint arXiv:2503.05379*.
- Zhao, J.; Yang, Q.; Peng, Y.; Bai, D.; Yao, S.; Sun, B.; Chen, X.; Fu, S.; Wei, X.; Bo, L.; et al. 2025b. HumanOmni: A Large Vision-Speech Language Model for Human-Centric Video Understanding. *CoRR*.
- Zhao, Z.; and Liu, Q. 2021. Former-dfer: Dynamic facial expression recognition transformer. In *Proceedings of the 29th ACM international conference on multimedia*, 1553–1561.

A Human Evaluation on Explanation Quality

Following the human evaluation described in Section 5.3, we detail here the two main purposes: (1) to assess the alignment between human judgments and a closed-source LLM (e.g., GPT-4.1-mini) regarding the representative emotion in each text explanation, and (2) to compare how well the explanations from different models capture the emotional content expressed in the video.

A.1 Survey Questions

We provided three questions for each sample, as show in Figure 5 and 4.

Question 1 asks for the representative emotion conveyed in the generated explanation. When tested on the output from our proposed model, the predictions showed 92.9% agreement with GPT-4.1-mini’s predictions. When considering the top 2 candidate emotions, the agreement increased to 96.4%.

Questions 2 and 3 evaluate how well the explanations reflect the emotion expressed in the video, not just by labeling an emotion. In other words, the focus is on how effectively the explanation interprets and describes the emotional signals embedded in the video. Specifically, each question is associated with one of the comparison models. These questions are rated on a 1-5 scale.

- 5 = Reflects the emotion in the video very well
- 4 = Reflects the emotion in the video well
- 3 = Somewhat reflects the emotion in the video
- 2 = Poorly reflects the emotion in the video
- 1 = Does not reflect the emotion in the video at all

Q1. {Emotion Explanation} What emotion does the explanation represent?

Q2&3. (Input Video) {Emotion Explanation} How well does this explanation reflect the emotion in the video?

Figure 4: Two questions used to evaluate the emotion conveyed in the generated explanation. {Emotion Explanation} can be a generated output from each R1-Omni and our model. (Input Video) means video is accompanied by the question.

13257

In the video, we see a female character with her eyes slightly closed and her mouth slightly parted. Her facial expression appears somewhat sad or contemplative. In the audio, it is impossible to determine the emotional state of the person. In the text, it is impossible to judge the emotional state based on the subtitle content.

What emotion does the explanation represent? *

☒ angry

☐ disgust


☐ happy

☐ sad

☐ surprise

☐ neutral

☐ fear



How well does this explanation reflect the emotion in the video? *

The video emotion analysis report shows that the female character in the video is in a dimly lit indoor environment. She wears a plaid shirt, with her hair neatly tied back in a ponytail, appearing somewhat reserved and contemplative. Her gaze shifts from time to time, revealing inner doubt or unease. Overall, she exhibits a complex emotional state mixed with confusion, worry, and slight sadness, suggesting she is facing challenges or setbacks in her life.

☐ 5 = Reflects the emotion in the video very well

☐ 4 = Reflects the emotion in the video well

☒ 3 = Somewhat reflects the emotion in the video

☐ 2 = Poorly reflects the emotion in the video

☐ 1 = Does not reflect the emotion in the video at all

How well does this explanation reflect the emotion in the video? *

In the video, we see a female character with her eyes slightly closed and her mouth slightly parted. Her facial expression appears somewhat sad or contemplative. In the audio, it is impossible to determine the emotional state of the person. In the text, it is impossible to judge the emotional state based on the subtitle content.

☐ 5 = Reflects the emotion in the video very well

☒ 4 = Reflects the emotion in the video well

☐ 3 = Somewhat reflects the emotion in the video

☐ 2 = Poorly reflects the emotion in the video

☐ 1 = Does not reflect the emotion in the video at all

Figure 5: The actual survey platform provided to participants in the human study.

Dataset	Angry	Happy	Sad	Surprise	Anxiety	Disgust	Fear	Neutral	Contempt	Disappointment	Helplessness
EMER	115	92	43	23	59	-	-	-	-	-	-
MERR-Fine	875	476	1,091	213	-	-	85	1,352	40	-	-
+ DFEW (Train)	2,834	1,534	1,566	1,558	2,769	146	1,195	515	38	582	737
+ MAFW (Train)	2,299	1,046	1,397	1,396	1,939	123	780	380	38	653	645
+ Complemented	-	-	-	-	-	600	-	500	800	200	200

Table 5: Statistics of the existing and augmented datasets used for training ERV. Datasets with a + in their names refer to augmented datasets consisting of emotional text description–emotion label pairs.

A.2 Statistics on the Human Assessment Scores

Table 6 presents the mean and standard deviation of the human assessment scores. The evaluated explanations had already been identified by the closed-source LLM as expressing the target emotion. In other words, these samples were included in our proposed EEA metric. Overall, these samples demonstrated performance above a score of 3, indicating that the explanations appropriately describe the emotions conveyed in the videos.

Model	Assessment Score
R1-Omni	3.42 ± 0.34
Ours	3.52 ± 0.41

Table 6: Mean and standard deviation of human assessment scores for emotional explanation quality across models.

B Augmenting Dataset for Training ERV

B.1 Details about Augmenting Emotional Text Descriptions

To address the limited dataset sizes of EMER and MERR-Fine, we augmented emotional text descriptions paired with emotion labels. Using the MAFW and DFEW training sets along with the R1-Omni model, we generated emotional descriptions for each video. However, as shown in Figure 1, the R1-Omni model fails to generate explanations that align with the ground-truth emotions. However, it is still capable of producing emotionally rich textual descriptions in and of themselves, even if they do not match the annotated labels.

Read the reasoning content and respond with the appropriate emotion in {Emotion List}. Primarily output a single dominant emotion. If necessary, output up to 2 dominant emotions, separated by commas. Reply with only the emotion word.

Reason: {Explanation} Answer Emotion:

Figure 6: Prompt used to generate the emotion conveyed in the generated explanation. {Emotion List} is a permutation of the GT emotion label from MAFW, and {Explanation} refers to the explanation generated by R1-Omni.

Leveraging this property, we generated pseudo-emotional text descriptions from the DFEW and MAFW datasets. We

then assigned emotion labels to these descriptions using GPT-4.1 with the prompt shown in Figure 6. This prompt is designed to output multi-label annotations from the given descriptions. The dataset statistics of MAFW and DFEW in Table 5 reflect overlapping multi-label annotations. As shown in the table, the resulting dataset exhibits class imbalance, with certain emotions such as *disgust*, *contempt*, *neutral*, *disappointment*, and *helplessness* being underrepresented.

To mitigate this issue, we generated additional emotional text descriptions specifically for these underrepresented classes using GPT-4.1-mini, guided by the prompt illustrated in Figure 7. Two-shot examples, randomly selected from the pseudo dataset for each specific emotion, are provided to the LLM to generate the "Complemented" dataset. This dataset complements class-imbalanced emotion descriptions.

You are an expert in audiovisual analysis with a talent for describing emotions. Goal: Carefully interpret the video’s visual, auditory, and contextual elements, then portray the Target Emotion in vivid language.

Example 1 (Target Emotion: {Specific Emotion}):
{Example Description 1}

Example 2 (Target Emotion: {Specific Emotion}):
{Example Description 2}

Now, based on the examples and the guidelines above, generate a new description that effectively conveys the following target emotion:

Target Emotion: {Specific Emotion}

Figure 7: Prompt used to generate emotion-oriented explanations from audiovisual content. {Emotion Explanation} refers to model outputs (e.g., R1-Omni or our model), and the (Input Video) refers to the input video used in the task.

E Training Details

We provide a more detailed description of the training configurations used in the SFT and GRPO training as shown in Table 8. We evaluated Emotion-LLaMA using the provided DFEW features and the checkpoint specialized for the emotion reasoning task. In the zero-shot evaluation on DFEW reported by Emotion-LLaMA, the model generated reasoning outputs for only about 40% of the test samples, so the reasoning metric could not be clearly measured.

configurations	SFT	SFT (7B)	GRPO	GRPO (7B)
Number of frames		8 uniformly sampled frames		
Image Resolution	224 × 224	384 × 384	224 × 224	384 × 384
Lr scheduler	cosine	cosine	cosine-decay	cosine-decay
Optimizer		AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.999$)		
DeepSpeed		ZeRO-3		
Training precision		bf16		
Total epochs	5	10	2	2
Warmup ratio	0.03	0.03	-	-
Learning rate	2e-5	2e-5	1e-6	1e-6
Batch size	32	32	16	16
Gradient accumulation	1	1	2	2
Trainable params	language model, vision tower, visual projector, audio projector, text gate module	language model, vision tower, visual projector, audio projector, text gate module	language model, visual projector, audio projector, text encoder	language model, visual projector, audio projector, text encoder
Vision Encoder	siglip-base (patch16-224)	siglip-so400m (patch14-384)	siglip-base (patch16-224)	siglip-so400m (patch14-384)
Audio Encoder		whisper-large-v3		

Table 8: Training specific configurations

F Computing Details of Explanation Reward

Algorithm 1 details the procedure for computing the Explanation Reward R_E for the output O generated by π_θ over G generations, as introduced in Figure 2 and Equation 4.

Algorithm 1: Pseudocode for Explanation Reward

Input:

G : Number of generations for GRPO training

$O = \{o_i\}_{i=1}^{N_{gen}}$: Output sequences from π_θ

ERV : ERV module returning logits

e_{gt} : Ground-truth emotion label

Parameters:

$\tau = 0.5$: Probability threshold for emotion selection

$k = 2$: Maximum number of emotions selected per sentence

Intermediate Variables:

E_i : Extracted explanation from o_i between `<think>` and `</think>`

$S_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,l}\}$: Set of sentences split from E_i

c_i : Number of sentences in S_i consistent with e_{gt}

$N_{i,neu}$: Number of neutral sentences in S_i

EL_k : Classified emotion label from $s_{i,k}$

$R_{i,E}$: Explanation Reward from o_i

Output:

R_E : List of all computed Explanation Rewards

```

1: Initialize  $R_E \leftarrow []$ 
2: for  $i = 1$  to  $G$  do
3:   Split  $E_i$  into sentences  $S_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,l}\}$ 
4:    $c_i \leftarrow 0, N_{i,neu} \leftarrow 0$ 
5:   for each  $s_{i,k}$  in  $S_i$  do
6:      $u_{i,k} \leftarrow \text{torch.sigmoid}(ERV(s_{i,k}))$ 
7:     if  $\sum(u_{i,k} > \tau) == 0$  then
8:        $EL_k \leftarrow \text{Top-1 label from } u_{i,k}$ 
9:     else if  $\sum(u_{i,k} > \tau) == 1$  then
10:       $EL_k \leftarrow \text{Label corresponding to the argmax of}$ 
11:         $u_{i,k}$ 
12:      else
13:         $EL_k \leftarrow \text{Top-}k \text{ labels from } u_{i,k}$ 
14:      end if
15:      if "neutral"  $\in EL_k$  then
16:         $N_{i,neu} \leftarrow N_{i,neu} + 1$ 
17:      end if
18:      if  $e_{gt} \in EL_k$  then
19:         $c_i \leftarrow c_i + 1$ 
20:      end if
21:    end for
22:    if  $e_{gt} == \text{"neutral"}$  then
23:       $R_{i,E} \leftarrow \frac{c_i}{|S|}$ 
24:    else
25:       $R_{i,E} \leftarrow \frac{c_i}{|S| - N_{i,neu}}$ 
26:    end if
27:     $R_E.append(R_{i,E})$ 
28: end for
29: return  $R_E$ 

```

G Qualitative Results

Figure 8, 9, 10, 11 show qualitative results comparing the generated explanations from R1-Omni and Ours. While R1-Omni produces the correct ground-truth emotion in the `<answer>` format, its predicted emotions from the closed-source LLM are often misaligned. In contrast, our model's explanations and final predicted answers demonstrate strong coherence.



Input video					
Input Speech	 "(Crying sound)"				
Baseline MLLM (R1-Omni)	<p><think>The video focuses on a middle-aged man wearing a white shirt, with slightly disheveled black hair. His expression shifts from initial confusion and contemplation to one of deep sadness and despair, eventually showing signs of sudden joy or relief. The entire process reveals the complex fluctuations in his emotions, as if he is experiencing an inner struggle and response to life's challenges. This emotional analysis not only highlights the authenticity and depth of his expression but also underscores the significance of emotional management and psychological resilience in overcoming difficulties.</think><answer>sad</answer></p>				<p>✗ Explanation: Happy</p> <p>✓ Answer: Sad</p>
Baseline MLLM With Emotional Rationale Verifier (Ours)	<p><think>In the video, we see a male character with his back to the camera. At the beginning of the video, he opens his eyes wide and looks down, as if crying or experiencing an emotional low. In the following scenes, he continues to cry, with tears streaming down his face, indicating that he is going through some sad things. In the audio, it seems like the character is crying because something bad has happened. In the text, the subtitle says, "I'm so sorry, I can't stand this anymore." This sentence expresses the emotional state of the male character. Based on the video clues of the male character opening his eyes wide and looking down, as well as the crying actions and tears streaming down his face in the video, it can be inferred that the male character is going through some sad things. Additionally, based on the audio clue describing the character crying due to something bad happening, it can be further confirmed that the male character's emotional state is one of sadness.</think><answer>sad</answer></p>				<p>✓ Explanation: Sad</p> <p>✓ Answer: Sad</p>

Figure 8: Generated emotional explanation comparison between R1-Omni vs. Ours about 15341.mp4 video in MAFW dataset



Input video					
Input Speech	 "(With a trembling voice)"				
Baseline MLLM (R1-Omni)	<p><think>In the video, a woman stands in front of a closed door, her facial expression tense and focused. Her brows are deeply furrowed, with a slight wrinkle of her nose, and her eyes reveal determination and anger. As she turns her head, her body language carries urgency and determination, as if she is facing an important decision or complaint. Her speech is fast-paced, with an increased volume and accelerated pace, indicating inner turmoil and a strong sense of dissatisfaction or anger. Overall, this woman experiences a significant emotional shift from tension to anger, reflecting intense internal conflict and emotional pressure.</think><answer>sad</answer></p>				<p>✗ Explanation: Angry</p> <p>✓ Answer: Sad</p>
Baseline MLLM With Emotional Rationale Verifier (Ours)	<p><think>In the video, we see a woman standing in an indoor environment. At the beginning of the video, her facial expression is relatively calm and focused, with slightly drooping corners of her mouth, indicating that she may be experiencing some emotional fluctuations or waiting for something. In the following scenes, her gaze shifts to the side, as if contemplating something. Her mouth moves up and down, suggesting that she is engaged in a conversation with someone. In the audio, it seems like the character is crying due to hearing bad news. In the text, the subtitle says, "I have to divorce my father because I'm pregnant." This sentence expresses the emotional state of the female character. Based on the visual clues of the woman's calm and focused facial expression, slight drooping corners of her mouth, and the contemplative and anticipatory behavior shown when looking at someone else, it can be inferred that she may be going through some emotional fluctuations or waiting for something. Additionally, based on the description of the character crying due to hearing bad news in the audio, it can be speculated that she may be going through a sad or worried emotion. Therefore, this sentence conveys the female character's emotional state of going through a difficult decision.</think><answer>sad</answer></p>				<p>✓ Explanation: Sad</p> <p>✓ Answer: Sad</p>

Figure 9: Generated emotional explanation comparison between R1-Omni vs. Ours about 11003.mp4 video in MAFW dataset

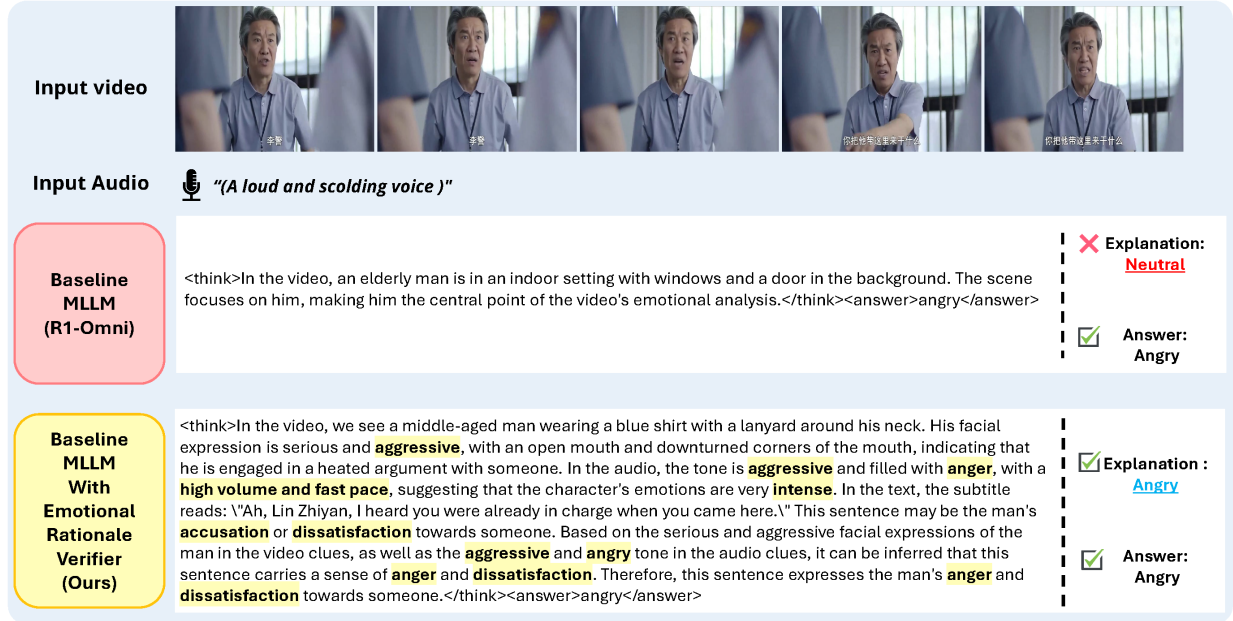


Figure 10: Generated emotional explanation comparison between R1-Omni vs. Ours about 10360.mp4 video in MAFW dataset

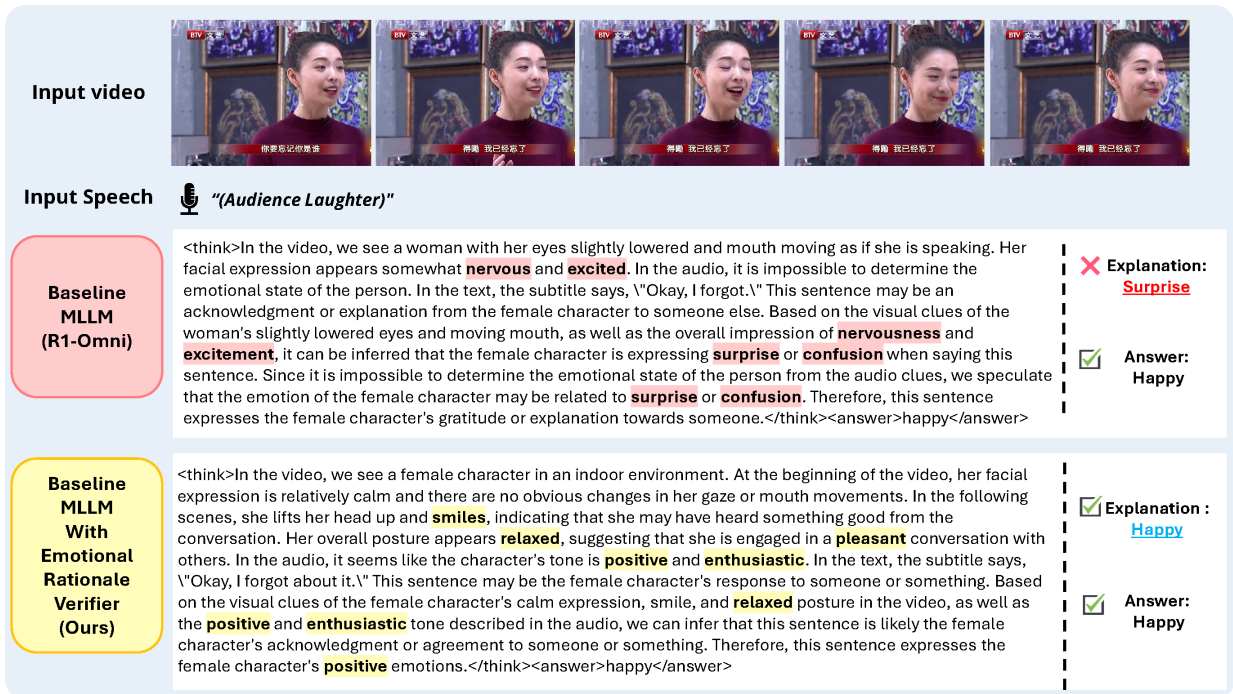


Figure 11: Generated emotional explanation comparison between R1-Omni vs. Ours about 15472.mp4 video in MAFW dataset