# Learning What to Attend First: Modality-Importance-Guided Reasoning for Reliable Multimodal Emotion Understanding

Hyeongseop Rha,* Jeong Hun Yeo,* Junil Won, Se Jin Park, and Yong Man Ro†

Integrated Vision and Language Lab, KAIST, South Korea

291 Daehak-ro, Yuseong-gu, Daejeon 34141, Republic of Korea

{ryool_1832, sedne246, dnjswnsdlf48, jinny960812, ymro}@kaist.ac.kr

## Abstract

*In this paper, we present Modality-Importance-Guided Reasoning (MIGR), a framework designed to improve the reliability of reasoning-based multimodal emotion understanding in multimodal large language models. Although existing methods have advanced emotion understanding, they often suffer from reasoning drift: models gradually rely on their own generated text instead of multimodal evidence, and their explanations are overly shaped by visually initiated reasoning paths. To address these issues, we introduce Modality Importance (MI), a simple yet effective mechanism for identifying the emotion-dominant modality. Using MI, MIGR reorganizes reasoning sequences so that explanations begin from the modality most critical to the target emotion, preventing early reasoning from being misled by less informative cues. Our two-stage framework—comprising modality-aligned supervised fine-tuning and modality-aware reward optimization—encourages models to generate emotionally grounded, causally relevant, and coherence-preserving explanations. Experimental results on the DFEW benchmark show that MIGR substantially improves reasoning reliability, decreasing instances of correct predictions accompanied by emotionally inconsistent explanations from 18.10% to 7.37%. These results confirm the benefit of initiating reasoning from the emotion-dominant modality.*

## 1. Introduction

As human–machine interaction becomes increasingly prevalent through AI assistants, smart speakers, and wearable devices, understanding human emotions has become crucial for achieving empathetic and socially aware communication. To address this need, Multimodal Emotion Recognition (MER) [1–8] has emerged as a key research di-
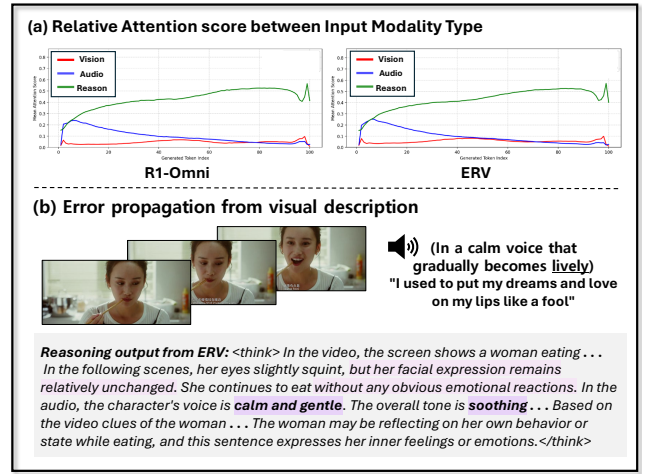


Figure 1. Failure pattern of the baseline multimodal reasoning model. (a) The model's attention gradually drifts from external multimodal inputs to previously generated text. (b) Due to the fixed training order: visual description → audio description → reasoning, the model often begins with a visual-based summary; when this initial step is inaccurate or emotionally irrelevant, subsequent reasoning becomes increasingly misaligned.

rection, leveraging complementary cues from visual, acoustic, and textual modalities to infer emotional states more accurately than unimodal approaches [9]. Recent progress in multimodal integration with Large Language Models (*i.e.*, MLLMs) has significantly improved MER performance, enabling models to more effectively leverage the complementary cues [10–14]. Building on this foundation, the introduction of Reinforcement Learning (RL) optimization has further extended MER from recognition to reasoning-based emotion understanding, allowing models to articulate why an emotion arises based on multimodal evidence [15].

Despite the recent progress, reasoning-based emotion understanding with MLLMs remains relatively underexplored. Our aim is to achieve reliable emotional reasoning, which we define as the ability to generate explanations that

---

*Equal contribution

†Corresponding author

are consistent with the model's final emotion prediction. It is an essential capability for emotional AI, which must not only predict emotions accurately but also justify them with consistent and trustworthy explanations in ambiguous multimodal scenarios. However, to the best of our knowledge, only one prior work has attempted to improve such reliability [16]. It introduces RL optimization that rewards MLLMs when an LLM judges the generated reasoning to be consistent with the prediction. While this approach shows some improvements, it remains an early step toward reliable emotional reasoning, and substantial further advances are still required.

In this context, we begin by exploring where this recent attempt [16] falls short. As shown in Figure 1(a), the model's attention gradually shifts from multimodal inputs (*i.e.*, audio and visual cues) to its own generated text as reasoning progresses. This shift causes the model to anchor its subsequent reasoning on whatever emotional cues appear in the initial generated text. Moreover, the reasoning process typically starts from generating a video description, because the training data are organized in a fixed order in which visual descriptions precede audio descriptions and the final reasoning. However, such a visual-first reasoning path may not reflect the true emotional driver, especially in cases where emotion is conveyed primarily through vocal tone or semantics. When the initial visual-based reasoning is irrelevant or inaccurate, the subsequent reasoning becomes increasingly misaligned, as illustrated in Figure 1(b). As a result, the initial step of reasoning becomes disproportionately influential. Similar issues appear in other domains, where early cues steer reasoning away from the underlying evidence [17].

To address the issue, we propose Modality-Importance-Guided Reasoning (MIGR), a method that aims to initiate reasoning from the emotion-dominant modality (*i.e.*, the modality most strongly associated with the target emotion) and maintain emotion-coherent reasoning throughout inference. To achieve this, we introduce Modality Importance (MI) estimation that identifies the modality most critical for accurate emotion recognition by comparing the model's predicted emotion under audio-only, visual-only, and audio–visual inputs. If the model accurately predicts the emotion under audio-only input but fails under visual-only or audio–visual conditions, this suggests that the visual modality may be irrelevant or even misleading. In this case, the audio modality is considered emotion-dominant.

Building on this MI and recent reasoning-based MLLM training paradigms [15, 16], we design our method as a two-stage learning framework consisting of (i) modality-aligned Supervised Fine-Tuning (SFT) and (ii) modality-aligned reward optimization. In the modality-aligned SFT stage, we leverage the MI to reorganize the original reasoning text into modality-specific segments and then reorder

according to emotion-dominant modality. Then, we train on the data so that it starts to generate emotion-relevant reasoning consistently at the beginning. In the modality-aligned reward optimization, we depart from conventional RL reward designs [15] that primarily focus on reasoning format or answer accuracy, and instead enforce modality-aware constraints through two complementary rewards. The modality-aligned order reward encourages the model to initiate its reasoning from the emotion-dominant modality, thereby reinforcing a modality-aligned reasoning structure. The modality-grounded reasoning reward further promotes generating emotion-dominant modality reasoning text that is causally and semantically aligned with the target emotion.

In this work, we make the following contributions: (i) We introduce the MI for identifying the emotion-dominant modality in multimodal emotion reasoning. It provides a simple yet effective way to diagnose misleading and informative modalities for modality-aligned emotion reasoning. (ii) We propose MIGR which restructures reasoning data so that the model initiates reasoning from the most informative modality. MIGR further reinforces this modality-aligned reasoning through modality-aware reward optimization, producing more stable and emotion-grounded explanations. (iii) MIGR achieves substantial improvements in reasoning reliability on the DEFW benchmark, including a 15.42% gain in Explanation–Prediction Consistency and a reduction of emotionally incorrect reasoning from 18.10% to 7.37%. These gains demonstrate the effectiveness of emotion-dominant modality alignment for robust multimodal emotional reasoning.

## 2. Related Work

### 2.1. Multimodal Emotion Recognition

#### 2.1.1. Traditional MER Approaches

MER aims to understand human affective states by analyzing multimodal signals that appear across diverse expressive scenarios. The development of MER was initially driven by the introduction of unimodal emotion-related datasets [1, 4, 7, 8, 18], which provided isolated affective cues such as audio, visual, or speech signals. Leveraging these datasets enabled researchers [19–23] to quantitatively analyze emotional expressions within each modality, leading to the development of modality-specific encoders that effectively model and encode unimodal emotional information. Subsequently, the introduction of multimodal emotion datasets [2, 3, 5, 6, 24] further accelerated the progress of MER by offering richer combinations of affective cues—including audio, visual, textual, and physiological signals. These datasets opened the door for modeling cross-modal interactions, driving the development of fusion strategies [25, 26] designed to integrate complemen-

tary information across modalities and improve recognition robustness.

### 2.1.2. Multimodal LLM-based MER

With the recent success of MLLMs across a wide range of domains, these models have also begun to be applied to emotion understanding tasks. The introduction of multimodal emotion-descriptive datasets such as EMER [27] has further accelerated progress in MLLM-based MER, enabling MLLMs to approach emotion recognition as a natural language–based reasoning and explanation task [10–14, 28]. However, such emotion descriptions require manual annotation, making large-scale expansion difficult. To address this limitation, recent studies [11, 12, 29] have proposed leveraging MLLMs to automatically generate emotion descriptions as pseudo-labels, enabling scalable construction of description-enriched datasets. Building upon these expanded datasets, several works [11, 29] have explored generating and interpreting emotion explanations grounded in video and audio cues. In parallel, architectural advances [28, 30] have further improved the accuracy of emotion recognition by enhancing multimodal representation learning and reasoning capabilities.

### 2.1.3. Reasoning with Multimodal LLMs

Despite significant progress in MLLM-based MER, these approaches still rely heavily on well-designed emotional annotations. This dependency has motivated growing interest in training paradigms that reduce reliance on explicit labels and instead strengthen reasoning through alternative supervision. In this context, RL-based post-training has recently emerged as an effective strategy for enhancing MLLMs' reasoning abilities while mitigating label dependence. In particular, GRPO [31] and Verifiable Reward [32] have demonstrated strong generalization not only in math and coding but also across a wide range of multimodal domains [15, 33–39]. However, many existing methods primarily focus on optimizing the final answer, while the fidelity and consistency of the intermediate reasoning process often receive insufficient attention, an issue highlighted by [40]. To address this, recent studies [41, 42] have introduced consistency-aware methodologies designed to better align reasoning with final outputs. Yet, such reasoning consistency problems remain largely unexplored in inherently ambiguous domains such as MER. Because emotional interpretation depends on subtle and context-dependent multimodal cues, reliable emotional reasoning becomes especially important. In particular, the generated explanations must faithfully support the model's final emotion prediction. Nevertheless, only one prior work has attempted to improve such reliability in MER [16], introducing an RL objective that rewards MLLMs when an LLM judge deems the generated reasoning consistent with the prediction. While this approach represents an important initial step, substantial further progress is required to achieve robust and trustworthy emotional reasoning.

## 3. Method

In this section, we introduce MIGR, a method that learns which modality to attend to first based on MI. Through this design, we aim to mitigate text biases in the initial reasoning stage and ensure that early reasoning remains grounded in the emotionally dominant modality. Our method consists of two learning stages. First, we present modality-aligned SFT, which enables the model to learn an initial reasoning direction by beginning from the modality identified as emotion-dominant by the MI label. Second, we introduce Modality-Aware Reward Optimization, which aims to reinforce modality prioritization and emotion-grounded reasoning while providing stability during optimization. To this end, we propose two rewards: modality-aligned order reward and modality-grounded reasoning reward. The following sections describe each training stage in detail.

### 3.1. Modality-Aligned SFT

Recent studies [15] have demonstrated the effectiveness of SFT as a cold-start in emotion understanding tasks, showing that even a small amount of high-quality supervision from the Explainable Multimodal Emotion Reasoning (EMER) dataset [27] can yield substantial performance gains. Motivated by these findings, we also initiate our training with a cold-start strategy.

#### 3.1.1. Emotion-Consistent Data Augmentation

Although EMER [27] provides high-quality supervision, the SFT stage can be further improved by incorporating samples with particularly clear and consistent emotional cues. To enhance early-stage supervision, we therefore augment our training set with additional emotion-consistent samples. To identify such samples, we extract Facial Action Units (FAUs) from the facial frames of each video, which provide strongly correlated emotional cues [43]. We then evaluate whether the extracted FAU patterns align with the target emotion based on a predefined FAU–emotion mapping table [11]. Only samples whose FAU patterns exactly match the target emotion are selected. For these filtered samples, we construct reasoning–answer pairs using FAU-consistent emotional evidence to guide the reasoning content.

#### 3.1.2. Modality Importance (MI)

We now introduce MI that determines the modality most strongly correlated with the target emotional signal. The key idea is to assess the contribution of each modality to the model's emotional understanding by comparing how accurately the model can infer the target emotion when different modality combinations (i.e., audio-only, visual-only,
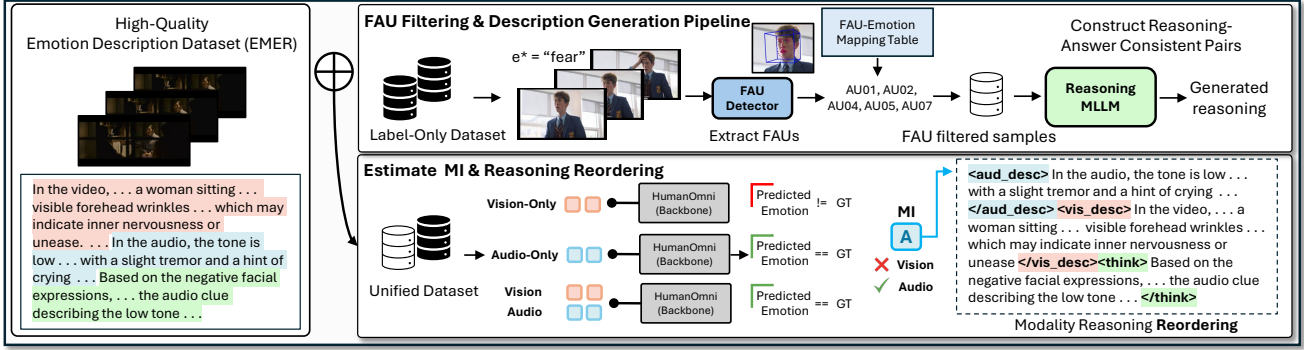
Figure 2. Overview of the proposed data construction pipeline, including FAU-based emotion-consistent data augmentation, MI estimation, and MI-guided modality-specific reasoning reordering.

and audio–visual) are used as inputs. By contrasting the model's behavior across these modality combinations, we identify which modality is the emotion-dominant in each audio–visual sample.

### 3.1.3. MI-Based Modality Reasoning Reordering

After constructing both the high-quality human-annotated dataset (EMER) and the FAU-filtered additional dataset, we merge them to form a unified training set. Using this combined data, we extract the MI for every sample. Once the emotion-dominant modality is determined, we reorganize the structure of the reasoning text, as shown in Figure 3.

Specifically, for each training example, we decompose the reasoning text into two modality-specific reasoning texts: an audio-based reasoning text and a video-based reasoning text. To explicitly mark these modality-specific parts, we employ two special tokens, <aud_desc> and <vis_desc>, placing the corresponding token at both the beginning and the end of each audio- and video-based text. We then reorder these text according to the MI. If audio is the emotion-dominant modality, the audio-based reasoning text is placed first, followed by the video-based text; conversely, if video is emotion-dominant, the video-based text is placed before the audio-based text. For samples where audio and video are equally informative, we include both orderings, *i.e.*, one with audio-first and one with video-first. By leveraging the reasoning text reorganized with the dominant modality during SFT, the model learns to initiate its reasoning process from the modality that carries the strongest emotional signal.

### 3.2. Modality-Aligned Reward Optimization

Following the cold-start stage, we refine the model's reasoning process through Modality-Aligned Reward Optimization. Prior GRPO-based approaches [15, 16] typically employ two reward types: an answer reward and a format reward. Since these rewards focus only on generating a correct answer and producing well-structured rea-

soning, they provide no explicit guidance on what kind of reasoning should be produced. Therefore, these two rewards are insufficient to preserve the emotion-dominant modality-prioritized reasoning learned during SFT. To reinforce MI-guided reasoning throughout optimization, we introduce two rewards: the modality-aligned order reward and the modality-grounded reasoning reward.

### 3.2.1. Modality-Aligned Order Reward

The modality-aligned order reward encourages the model to generate reasoning in a modality order consistent with the MI $m$. During optimization, if audio is the emotion-dominant modality, the reward is granted when the model initiates its reasoning with the <aud_desc> token and the target emotion is detected within the <aud_desc> token. Similarly, if the video is emotion-dominant, the reward is granted when the reasoning begins with the <vis_desc> token and the target emotion appears within the <vis_desc> token. We denote the modality-specific token corresponding to the MI $m$ as <M>. For this token, we define $E_m(y)$ as the set of predicted emotion labels obtained by applying an emotion classifier $e_{classifier}$ to its reasoning sentences. Under this definition, the modality-aligned order reward is given by:

$$R_{\text{MAO}}(y, m, e^*) = \begin{cases} 1, & \text{if } y \text{ starts with <M>}, e^* \in E_m(y) \\ 0, & \text{otherwise.} \end{cases}$$

(1)

### 3.2.2. Modality-Grounded Reasoning Reward

The modality-grounded reasoning reward evaluates the semantic alignment between each modality-specific reasoning text and the target emotion. For each sample, we measure how well the audio-based and video-based reasoning texts provide emotion-consistent evidence relative to the ground-truth emotion. The reward is then computed based on the MI: if audio is the emotion-dominant modality, we

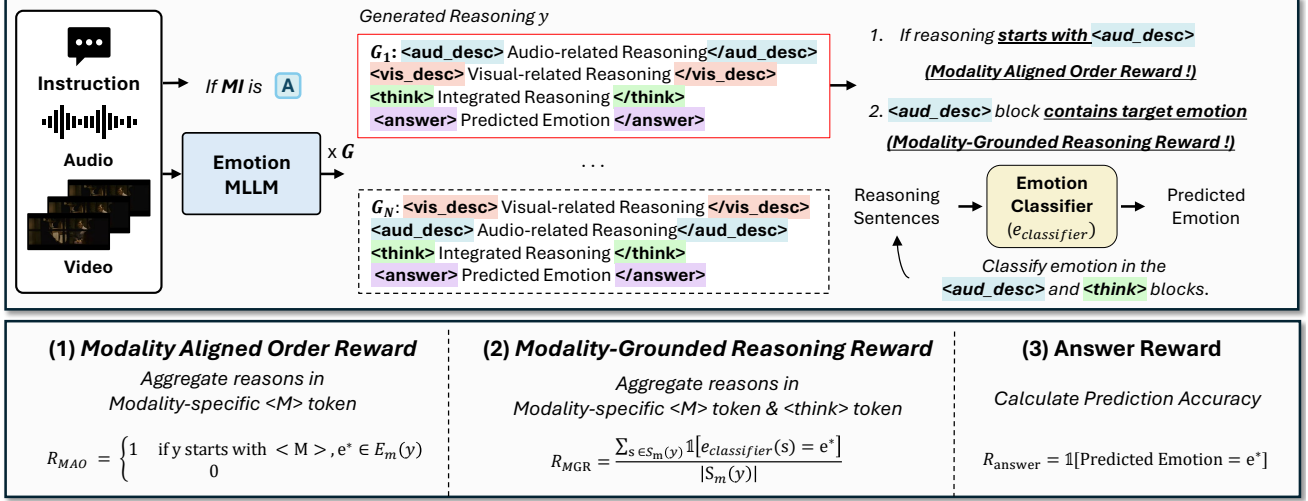**MIGR Reward Computation in the Training pipeline**



Figure 3. Illustration of the three rewards in MIGR: the Modality-Aligned Order Reward enforces MI-consistent reasoning order; the Modality-Grounded Reasoning Reward ensures emotion-consistent modality-specific reasoning; and the answer reward guarantees correct final emotion prediction.

aggregate the reasoning sentences from the `<aud_desc>` and `<think>` tokens. Likewise, if video is emotion-dominant, we aggregate the reasoning sentences from the `<vis_desc>` and `<think>` tokens. An emotion classifier $e_{classifier}$ is then applied to this aggregated set of reasoning sentences, denoted as $S_m(y)$, to obtain the predicted emotion labels. The correction ratio computed from these predictions is used as the reward. This reward function can be formulated as follows:

$$R_{\text{MGR}}(y, m, e^*) = \frac{\sum_{s \in S_m(y)} \mathbb{1}\left[e_{\text{classifier}}(s) = e^*\right]}{|S_m(y)|} \quad (2)$$

In addition to these two modality-focused rewards, we further include an answer reward [15] ($R_{\text{answer}}$) to ensure the correctness of the final emotion prediction. In total, MIGR employs three rewards during optimization as follows:

$$R_{\text{total}} = R_{\text{MAO}} + R_{\text{MGR}} + R_{\text{answer}}. \quad (3)$$

As shown in Figure 3, MIGR integrates these three rewards to guide modality-aligned reasoning during optimization.

## 4. Experimental Setup

### 4.1. Datasets

We utilize three multimodal emotion dataset to train our model and validate its effectiveness. **EMER** [27] provides high-quality human-verified reasoning–label pairs and is used for cold-start SFT. **DFEW** [7] and **MAFW** [6] are large-scale benchmarks covering diverse facial expressions

in video, audio, and text modalities. A detailed description of dataset composition and statistics is provided in the Appendix 7.1.

### 4.2. Evaluation Metrics

Following prior work [11, 16], we report UAR and WAR for recognition accuracy, and three complementary metrics—EEA, EPC, and FCR—to quantify the emotional coherence of generated reasoning. These metrics jointly assess how well the explanation aligns with the target emotion and the model's prediction. Formal definitions and annotation protocols are included in the Appendix 7.2.

### 4.3. Implementation Details

**Pre-processing.** For video input, we resize each frame to $384 \times 384$ and uniformly sample 8 frames per clip. For audio, we use 16 kHz waveforms and convert them into 128-channel mel-spectrograms using a window size of 25 ms and a hop size of 10 ms.

**Architecture.** Following the design of R1-Omni [15], our framework consists of SigLIP [45] as the vision encoder, Whisper-large-v3 [46] as the audio encoder, and BERT [47] as the text encoder. Each modality output is projected into the LLM embedding space through its corresponding visual and audio projectors, implemented as two linear layers to match the dimensionality of the LLM representation. For reference, a 1-second audio clip yields approximately 50 audio features, while a single video frame produces 729 visual features. The LLM backbone is Qwen2.5-7B [48], and all model weights are initialized from the HumanOmni [30] checkpoint to preserve pretrained multimodal representations.

Table 1. Comparison with state-of-the-art methods on the MAFW and DFEW datasets. We report consistency metrics (FCR, EEA, EPC) and recognition accuracy (UAR, WAR). A dagger (†) denotes our re-implementation with identical training data for fair comparison.

| Model | MLLM Based | MLLM Params | MAFW | | | | | DFEW | | | | |
| | | | Consistency (%) ↑ | | | Accuracy (%) ↑ | | Consistency (%) ↑ | | | Accuracy (%) ↑ | |
| | | | FCR | EEA | EPC | UAR | WAR | FCR | EEA | EPC | UAR | WAR |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| *Non-Reasoning Models* | | | | | | | | | | | | |
| Former-DFER [21] | ✗ | - | – | – | – | – | – | – | – | – | 53.69 | 65.70 |
| MAE-DFER [22] | ✗ | - | – | – | – | 41.62 | 54.31 | – | – | – | 63.41 | 74.43 |
| S2D [44] | ✗ | - | – | – | – | 41.86 | 57.37 | – | – | – | 61.82 | 76.03 |
| Emotion-LLaMA [11] | ✓ | 7B | – | – | – | – | – | – | – | – | 64.21 | 77.06 |
| FaVChat [14] | ✓ | 7B | – | – | – | 49.07 | 60.97 | – | – | – | 68.17 | 77.88 |
| Omni-Emotion [12] | ✓ | 7B | – | – | – | 53.81 | 64.23 | – | – | – | 68.80 | 78.35 |
| HumanOmni [30] | ✓ | 7B | – | – | – | 52.94 | 68.40 | – | – | – | 74.86 | 82.48 |
| *Reasoning-Generating Models* | | | | | | | | | | | | |
| R1-Omni† [15] | ✓ | 7B | 40.22 | 44.97 | 55.19 | 47.27 | 65.36 | 52.01 | 55.39 | 60.70 | 69.69 | 76.67 |
| ERV [16] | ✓ | 7B | 50.98 | 54.70 | 73.06 | 47.01 | 65.19 | 62.06 | 65.50 | 73.53 | 68.88 | 75.81 |
| **MIGR (Ours)** | ✓ | 7B | 55.30 | 57.65 | 84.37 | 44.82 | 62.46 | 68.48 | 70.06 | 88.95 | 67.78 | 73.93 |

**Training and Evaluation.** MIGR consists of two training stages: the SFT stage and the GRPO stage. In the SFT stage, we construct an additional emotional reasoning dataset using the MAFW and DFEW training datasets. Based on the ERV [16], we align samples where the emotion represented by the reasoning output matches the target emotion. AU (Action Unit) information is extracted using the OpenFace toolkit, and AU sets are organized following the dataset construction pipeline of Emotion-LLaMA [11]. Using the Emotion–AU Table for emotion alignment verification, we incorporate 184 samples from MAFW and 253 samples from the DFEW training set as additional reasoning data. Training is conducted for 5 epochs with a cosine scheduler, a warmup ratio of 0.03, a learning rate of 2e-5, and a batch size of 32, using $8 \times$ NVIDIA A100 GPUs. In the GRPO stage, the gradient accumulation step is set to 2, the local batch size to 1, and the generation number (G) to 4, with a learning rate of 1e-6, and a cosine decay scheduler. The training is performed for 1 epoch, and evaluation is conducted with a temperature of 0.3.

## 5. Experimental Results

### 5.1. Main Result

#### 5.1.1. Comparison with the state-of-the-art methods

To assess the effectiveness of the proposed framework, we compare MIGR with existing emotion understanding methods on the MAFW and DFEW datasets (Table 1). We categorize previous methods into two groups: (i) *Non-Reasoning Models* that directly predict emotion labels, and (ii) *Reasoning-Generating Models* that produce textual explanations before making a prediction.

We first compare MIGR with other reasoning-generating models in terms of consistency metrics to validate that the reasoning texts generated by our model are emotionally coherent. As shown in Table 1, MIGR consistently outperforms ERV on all three consistency measures across

both datasets. On MAFW, MIGR improves over ERV by +4.32, +2.95, and +11.31 points in FCR, EEA, and EPC, respectively, achieving 55.30% FCR, 57.65% EEA, and 84.37% EPC. Similarly, on DFEW, MIGR attains 68.48% FCR, 70.06% EEA, and 88.95% EPC, surpassing ERV by +6.42, +4.56, and +15.42 points. We then compare recognition accuracy between non-reasoning models and reasoning-generating models. Among non-reasoning approaches, HumanOmni achieves the strongest performance with 68.40% WAR on MAFW and 82.48% WAR on DFEW, establishing a strong accuracy-oriented baseline. In contrast, all reasoning-generating models, including MIGR and ERV, exhibit slightly lower recognition accuracy than HumanOmni, suggesting that the explicit reasoning generation process can weaken the direct answer-prediction capability. When comparing MIGR with reasoning models, our framework attains 62.46% WAR on MAFW and 73.93% WAR on DFEW, which is somewhat lower than the best reasoning counterparts, but accompanied by substantially higher consistency scores.

#### 5.1.2. Analysis of Reasoning-Answer Inconsistency

To validate whether the improvement in emotional coherence of MIGR stems from more faithful reasoning, we further examine cases where the reasoning explanation contradicts the target emotion while the predicted answer remains correct. To this end, we utilize two indicators: (i) the proportion of all samples where reasoning is incorrect but the answer is correct (R× / A✓), and (ii) the proportion of such reasoning errors among correctly predicted samples.

As shown in the Table 2, MIGR reduces reasoning–answer inconsistency compared to reasoning-generating models, R1-Omni, and ERV across both datasets. On the DFEW dataset, the percentage of inconsistent cases drops from 24.66% (R1-Omni) and 13.75% (ERV) to 5.45% with MIGR, while the relative inconsistency among correct predictions decreases from 32.2% to

Table 2. Reasoning–answer inconsistency analysis on the MAFW and DFEW datasets. R×/A✓ represents cases where the reasoning is incorrect while the predicted emotion matches the ground truth.

| Model | R× / A✓ (%)↓ | | | |
| --- | --- | --- | --- | --- |
| | DFEW | | MAFW | |
| | All Samples | Among Correct | All Samples | Among Correct |
| R1-Omni [15] | 24.66 | 32.20 | 25.14 | 38.50 |
| ERV [16] | 13.75 | 18.10 | 14.21 | 21.80 |
| **MIGR (Ours)** | **5.45** | **7.37** | **7.16** | **11.50** |

Table 3. Emotion-wise accuracy comparison on the DFEW dataset. We report per-emotion recognition accuracy along with overall UAR and WAR. A dagger (†) denotes our re-implementation using identical training data.

| Method | Hap | Sad | Neu | Ang | Sur | Dis | Fea | UAR | WAR |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Emotion-LLaMA [11] | 93.05 | 79.42 | 72.47 | 84.14 | 72.79 | 3.45 | 44.20 | 64.21 | 77.06 |
| R1-Omni† [15] | 90.59 | 78.04 | 62.29 | 85.71 | 73.04 | 27.59 | 70.56 | 69.69 | 76.67 |
| ERV [16] | 88.14 | 76.72 | 62.85 | 84.56 | 78.84 | 31.03 | 60.00 | 68.88 | 75.81 |
| **MIGR (Ours)** | 86.50 | 79.10 | 60.04 | 82.03 | 66.89 | 31.03 | 68.89 | 67.78 | 73.93 |

7.37%. A similar trend is observed on MAFW, with reductions from 25.14% to 7.16% (overall) and 38.5% to 11.5% (among correct samples). These findings suggest that MIGR not only generates explanations that are more consistent with predicted emotions but also effectively mitigates spurious reasoning that leads to correct answers for the wrong reasons.

### 5.1.3. Emotion-wise Performance Analysis

To investigate the underlying cause of the reduced accuracy observed in Table 1 compared to existing methods, we conduct an emotion-wise performance analysis on the DFEW dataset—where Hap (Happiness), Sad (Sadness), Neu (Neutral), Ang (Anger), Sur (Surprise), Dis (Disgust), and Fea (Fear) are the seven evaluation categories—as presented in Table 3.

In this experiment, we identify two key findings. First, both reasoning-generating models, ERV and MIGR demonstrate substantially stronger performance than the non-reasoning model (Emotion-LLaMA [11]) on the challenging Disgust and Fear categories. This suggests that incorporating explicit reasoning allows models to better capture subtle or ambiguous emotional cues, particularly when predicting emotions such as Fear and Surprise. Second, although MIGR achieves accuracy levels comparable to other reasoning-based approaches and maintains stable performance across most of the remaining emotion classes, it exhibits a pronounced degradation on the Surprise category compared with ERV, with accuracy dropping from 78.84% to 66.89%. To more precisely understand this degradation on Surprise, we further analyze the prediction errors of MIGR in comparison with ERV, focusing on the confusion patterns for Surprise and identifying which emotion

Table 4. Analysis of the 47 DFEW samples where ERV correctly predicts Surprise but MIGR does not. The table shows MIGR's misclassification distribution and the frequency with which Surprise appears among the Top-2 reasoning-level predictions.

| Samples where ERV is correct (Surprise) but MIGR is incorrect ($N = 47$) | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| MIGR predicted label distribution | | | | | | | Surprise Recall (Top-2) |
| Predicted label | Hap | Neu | Ang | Dis | Fea | Unk | Sur |
| # samples | 2 | 18 | 3 | 1 | 21 | 2 | 39 |
| Ratio (%) | 4.3 | 38.3 | 6.4 | 2.1 | 44.7 | 4.3 | 83.0 |

Table 5. Ablation study of the MIGR framework on the DFEW dataset. Stage 1 evaluates the effect of augmentation and reasoning reordering during SFT, and Stage 2 examines the impact of the two MI-guided GRPO rewards.

| Method | Consistency (%)↑ | | | Accuracy (%)↑ | |
| --- | --- | --- | --- | --- | --- |
| | FCR | EEA | EPC | UAR | WAR |
| Stage 1: Supervised Fine-Tuning (SFT) | | | | | |
| **Baseline (Used-only EMER)** | 42.21 | 44.86 | 57.41 | 39.16 | 48.33 |
| + Data Augmentation | 58.99 | 59.97 | 85.19 | 54.60 | 62.80 |
| + Training Data Reordering | 59.72 | 61.13 | 84.80 | 56.04 | 64.77 |
| Stage 2: SFT followed by GRPO (init: Stage 1 final model) | | | | | |
| + Modality-Aligned Order Reward | 64.64 | 66.14 | 82.92 | 68.18 | 74.91 |
| + Modality-Grounded Reasoning Reward | 67.25 | 68.28 | 88.27 | 67.24 | 73.84 |
| + Both Rewards (MIGR) | 68.48 | 70.06 | 88.95 | 67.78 | 73.93 |

categories it is most frequently misclassified into.

### 5.1.4. Analyzing Performance Degradation on the Surprise Category

To analyze why Surprise shows particularly lower performance, we examine 47 samples on which ERV correctly predicts Surprise but MIGR fails. As shown in Table 4, MIGR predominantly misclassifies these samples as Neutral or Fear. This pattern aligns with the well-known perceptual similarity between Fear and Surprise, which often share facial cues such as widened eyes or raised eyebrows, making them challenging to distinguish even for humans. To further examine whether MIGR truly "loses" the Surprise signal, we analyze the generated reasoning using a closed-source LLM. Interestingly, although MIGR's final predictions are incorrect, Surprise still appears within the Top-2 reasoning-level emotion predictions for 83% of the samples. This indicates that MIGR retains Surprise-related evidence in its reasoning but ultimately assigns the final label to a neighboring emotion category. Overall, these findings suggest that MIGR's reduced Surprise accuracy may stem from the inherent ambiguity between Surprise and Fear, indicating that this aspect has room for further improvement.

### 5.2. Ablation study

To better understand the contribution of each component in MIGR, we conduct a stage-wise ablation analysis on the DFEW dataset, as summarized in Table 5. Our ablation follows the two-stage training pipeline described in Section 3 and evaluates how successive additions in both SFT and GRPO stages influence emotional consistency and classi-

**MIGR (Ours):** "<vis_desc> In the video, we see ... </vis_desc><aud_desc>In the audio, it is **impossible** to determine the character's emotional state. In the text, it is **impossible** to judge the emotional state based on the subtitle content.</aud_desc><think>The visual cues strongly indicate that ... . **Although the audio and subtitle content do not provide clear emotional signals, the unmistakable visual expression carries significant weight in emotion prediction. When vocal tone or verbal content are neutral or ambiguous, facial expressions often offer the most direct window into one's emotional state.** Here, the man's startled facial appearance combined with his attentive posture convincingly supports the conclusion that he is surprised by what he has heard.</think><answer> surprise </answer>"

**ERV–R1:** "<think>In the video, we see a male character standing in a dimly lit indoor environment. His facial expression is somewhat negative, his eyes are wide open, and he seems to be staring directly at someone or something outside the frame, giving off an overall sense of tension. In the audio, it is **impossible** to determine the character's emotional state. In the text, it is **impossible** to judge the emotional state based on the subtitle content.</think><answer>fear</answer>"

**R1-Omni:** "<think>In the video, we see a male character standing in a dimly lit indoor environment. His facial expression is relatively calm, but his eyes seem slightly confused, as if he is contemplating something. In the audio, there are no valuable emotional clues provided. In the text, it is **impossible** to determine the emotional state based on the subtitle content.</think><answer>surprise</answer>"

**MIGR (Ours):** "<aud_desc>In the audio, the character\'s voice is accompanied by obvious **sobbing.** ... </aud_desc> <vis_desc>In the video, ... her eyes are wide open and her facial expression is very negative, ... , indicating that she may be experiencing an unpleasant emotional state. ... </vis_desc><think> Based on the negative facial expression and wide-open eye sockets of the female character in the video clues, as well as the obvious **sobbing** and low tone in the audio clues, we can infer that this sentence may carry a sense of worry or nervousness. ... Therefore, this sentence expresses her expectations and hopes for the other person\'s understanding and support.</think><answer>sad</answer>"

**ERV–R1:** "<think>In the video, ... expressing extreme emotional distress. In the audio, the tone is urgent and the character\'s emotions are heightened. Combined with the text content, it can be felt that the character is very **angry and anxious**. Based on the woman\'s facial expressions becoming more negative, ... At the same time, based on the urgent tone and heightened emotions in the audio clues, ... </think><answer>angry</answer>"

**R1-Omni:** " <think>In the video, ... speaking forcefully or shouting, ... In the audio, the tone is urgent and the character\'s emotions are heightened. Combined with the text content, we can sense the character\'s **impatience and anxiety.** The subtitle in the text says, ... Based on the changes in the woman\'s facial expression ... it can be inferred that her emotions are becoming more **agitated and possibly negative**. At the same time, based on the description of the urgent tone and heightened emotions in the audio clues, as well as the expression of **impatience and anxiety** in the subtitle, ... it can be speculated that the woman is expressing her **concern and anxiety** about the person she is talking to.</think><answer>sad</answer>"
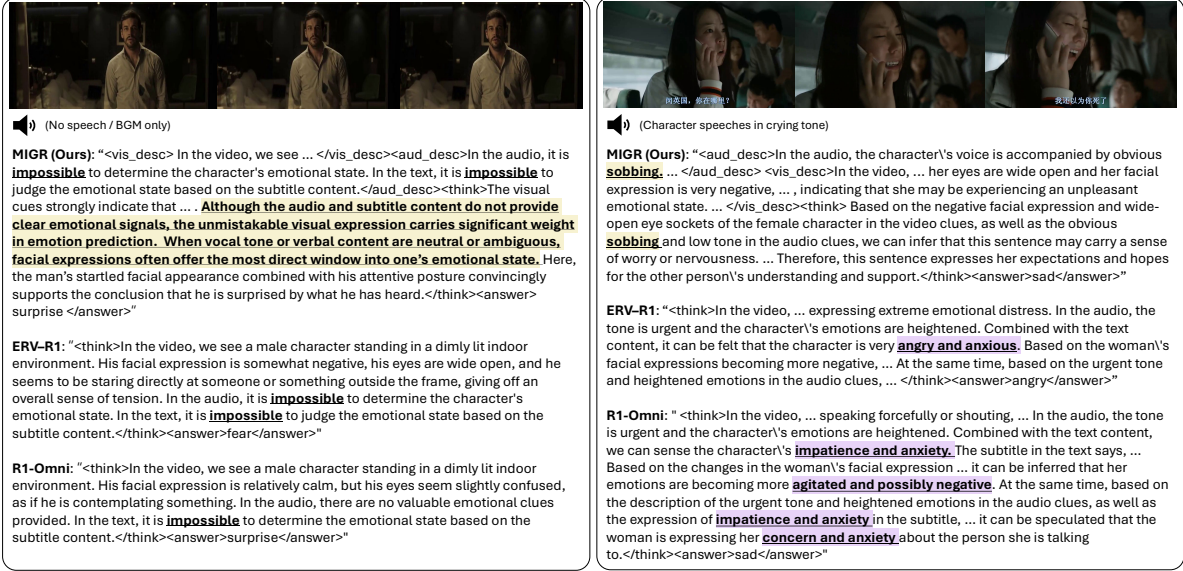
Figure 4. Qualitative comparisons of emotion reasoning. (Left) For a speechless sample, baseline models fail to infer the emotion due to missing audio cues, whereas MIGR correctly focuses on the visual modality and integrates it within the final <think> step to produce a coherent conclusion. (Right) In an audio-dominant sample, MIGR first identifies the key audio cue ("sobbing") and leverages it to interpret ambiguous visual information, leading to an accurate prediction. In contrast, other models misinterpret the visual cue (frowning) as anger and produce incorrect reasoning.

fication accuracy.

We begin with a baseline model trained solely on the EMER dataset, which provides high-quality but limited supervision. Introducing emotion-consistent data augmentation notably improves all consistency metrics, confirming the importance of supplying emotionally reliable samples during early-stage training. Subsequently, applying MI-based reasoning reordering yields additional gains by guiding the model to initiate its reasoning from the emotion-dominant modality. This validates our hypothesis that organizing reasoning sequences around MI helps mitigate visually biased or text-driven drift in the initial reasoning stage.

In the second stage, we examine how each proposed reward contributes to reasoning stability and coherence. The Modality-Aligned Order Reward encourages the model to begin reasoning from the modality prioritized by MI, restoring the modality-first structure learned during SFT. The Modality-Grounded Reasoning Reward further strengthens the emotional validity of modality-specific reasoning, promoting causal and semantically consistent evidence aligned with the target emotion. When both rewards are applied together, we observe the highest consistency across all metrics, demonstrating that the two rewards provide complementary benefits.

### 5.3. Qualitative Results

Figure 4 presents qualitative comparisons illustrating how MIGR utilizes modality-aligned reasoning to make more reliable emotion predictions. In the first example, the input contains no speech, making audio-based evidence unavailable. Baseline models, ERV, and R1-Omni, therefore fail to describe the emotional state because they are confused by audio or text cues. In contrast, MIGR correctly recognizes that the visual modality is dominant and structures its reasoning accordingly: the model begins with $<vis\_desc>$, identifies key facial cues, and then consolidates the interpretation in the $<think>$ step, leading to a coherent and accurate conclusion. The second example highlights the benefit of modality-grounded reasoning in audio-dominant scenarios. MIGR initiates its reasoning with the audio modality and identifies sobbing, a strong indicator of sadness. This early recognition enables the model to reinterpret the negative visual cues in the context of the audio evidence, resulting in the correct prediction. Baseline models, however, over-rely on visual information and incorrectly interpret the frowning expression as anger, leading to faulty reasoning.

## 6. Conclusion

In this work, we aimed to improve the reliability of reasoning-based multimodal emotion understanding by addressing the early-stage reasoning drift commonly observed in MLLMs. To this end, we introduced MI and proposed MIGR, a training framework designed to initiate and maintain reasoning from the modality most relevant to the target emotion. Our two-stage approach ensures that both the supervised and reinforcement learning processes are aligned with the emotion-dominant modality, leading to ex-

planations that are more coherent, emotionally grounded, and causally meaningful. Experiments on the MAFW and DFEW benchmark show that MIGR greatly reduces emotionally inconsistent reasoning. Overall, MIGR moves toward more trustworthy and interpretable multimodal emotional reasoning, highlighting the value of incorporating modality dominance into both data organization and optimization.

# References

[1] Jean Kossaifi, Georgios Tzimiropoulos, Sinisa Todorovic, and Maja Pantic. Afew-va database for valence and arousal estimation in-the-wild. *Image and Vision Computing*, 65:23–36, 2017. 1, 2

[2] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 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)*, pages 2236–2246, 2018. 2

[3] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations. *arXiv preprint arXiv:1810.02508*, 2018. 2

[4] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42:335–359, 2008. 2

[5] Olga Perepelkina, Evdokia Kazimirova, and Maria Konstantinova. Ramas: Russian multimodal corpus of dyadic interaction for affective computing. In *International Conference on Speech and Computer*, pages 501–510. Springer, 2018. 2

[6] Yuanyuan Liu, Wei Dai, Chuanxu Feng, Wenbin Wang, Guanghao Yin, Jiabei Zeng, and Shiguang Shan. 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*, pages 24–32, 2022. 2, 5, 12

[7] Xingxun Jiang, Yuan Zong, Wenming Zheng, Chuangao Tang, Wanchuang Xia, Cheng Lu, and Jiateng Liu. Dfew: A large-scale database for recognizing dynamic facial expressions in the wild. In *Proceedings of the 28th ACM international conference on multimedia*, pages 2881–2889, 2020. 2, 5, 12

[8] Steven R Livingstone and Frank A Russo. 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, 2018. 1, 2

[9] Yuntao Shou, Tao Meng, Wei Ai, and Keqin Li. Multimodal large language models meet multimodal emotion recognition and reasoning: A survey. *arXiv preprint arXiv:2509.24322*, 2025. 1

[10] Bohao Xing, Zitong Yu, Xin Liu, Kaishen Yuan, Qilang Ye, Weicheng Xie, Huanjing Yue, Jingyu Yang, and Heikki Kälviäinen. Emo-llama: Enhancing facial emotion understanding with instruction tuning. *arXiv preprint arXiv:2408.11424*, 2024. 1, 3

[11] Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and Alexander Hauptmann. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning. *Advances in Neural Information Processing Systems*, 37:110805–110853, 2024. 3, 5, 6, 7, 12

[12] Qize Yang, Detao Bai, Yi-Xing Peng, and Xihan Wei. Omniemotion: Extending video mllm with detailed face and audio modeling for multimodal emotion analysis. *arXiv preprint arXiv:2501.09502*, 2025. 3, 6

[13] Zhuozhao Hu, Kaishen Yuan, Xin Liu, Zitong Yu, Yuan Zong, Jingang Shi, Huanjing Yue, and Jingyu Yang. Feallm: Advancing facial emotion analysis in multimodal large language models with emotional synergy and reasoning. *arXiv preprint arXiv:2505.13419*, 2025.

[14] Fufangchen Zhao, Ming Li, Linrui Xu, Wenhao Jiang, Jian Gao, and Danfeng Yan. Favchat: Unlocking fine-grained facial video understanding with multimodal large language models. *arXiv preprint arXiv:2503.09158*, 2025. 1, 3, 6, 12

[15] Jiaxing Zhao, Xihan Wei, and Liefeng Bo. R1-omni: Explainable omni-multimodal emotion recognition with reinforcement learning. *arXiv preprint arXiv:2503.05379*, 2025. 1, 2, 3, 4, 5, 6, 7

[16] Hyeongseop Rha, Jeong Hun Yeo, Yeonju Kim, and Yong Man Ro. Emotion-coherent reasoning for multimodal llms via emotional rationale verifier. *arXiv preprint arXiv:2510.23506*, 2025. 2, 3, 4, 5, 6, 7, 12

[17] Mi Luo, Zihui Xue, Alex Dimakis, and Kristen Grauman. When thinking drifts: Evidential grounding for robust video reasoning. *arXiv preprint arXiv:2510.06077*, 2025. 2

[18] Yan Wang, Yixuan Sun, Yiwen Huang, Zhongying Liu, Shuyong Gao, Wei Zhang, Weifeng Ge, and Wenqiang Zhang. 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*, pages 20922–20931, 2022. 2

[19] Shiqing Zhang, Shiliang Zhang, Tiejun Huang, and Wen Gao. Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching. *IEEE transactions on multimedia*, 20(6):1576–1590, 2017. 2

[20] Shiqing Zhang, Xin Tao, Yuelong Chuang, and Xiaoming Zhao. Learning deep multimodal affective features for spontaneous speech emotion recognition. *Speech Communication*, 127:73–81, 2021.

[21] Zengqun Zhao and Qingshan Liu. Former-dfer: Dynamic facial expression recognition transformer. In *Proceedings of the 29th ACM international conference on multimedia*, pages 1553–1561, 2021. 6

[22] Licai Sun, Zheng Lian, Bin Liu, and Jianhua Tao. Maedfer: Efficient masked autoencoder for self-supervised dynamic facial expression recognition. In *Proceedings of the*

*31st ACM International Conference on Multimedia*, pages 6110–6121, 2023. 6

[23] Kateryna Chumachenko, Alexandros Iosifidis, and Moncef Gabbouj. 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*, pages 4673–4682, 2024. 2

[24] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and Kwanghoon Sohn. Context-aware emotion recognition networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10143–10152, 2019. 2

[25] Kartik Narayan, Vibashan VS, Rama Chellappa, and Vishal M Patel. Facexformer: A unified transformer for facial analysis. *arXiv preprint arXiv:2403.12960*, 2024. 2

[26] Licai Sun, Zheng Lian, Bin Liu, and Jianhua Tao. Hicmae: Hierarchical contrastive masked autoencoder for self-supervised audio-visual emotion recognition. *Information Fusion*, 108:102382, 2024. 2

[27] Zheng Lian, Licai Sun, Mingyu Xu, Haiyang Sun, Ke Xu, Zhuofan Wen, Shun Chen, Bin Liu, and Jianhua Tao. Explainable multimodal emotion reasoning. *CoRR*, 2023. 3, 5, 12

[28] Zheng Lian, Haoyu Chen, Lan Chen, Haiyang Sun, Licai Sun, Yong Ren, Zebang Cheng, Bin Liu, Rui Liu, Xiaojiang Peng, et al. Affectgpt: A new dataset, model, and benchmark for emotion understanding with multimodal large language models. In *Forty-second International Conference on Machine Learning*. 3

[29] Zheng Lian, Haiyang Sun, Licai Sun, Lan Chen, Haoyu Chen, Hao Gu, Zhuofan Wen, Shun Chen, Zhang Siyuan, Hailiang Yao, et al. Open-vocabulary multimodal emotion recognition: Dataset, metric, and benchmark. 2024. 3

[30] Jiaxing Zhao, Qize Yang, Yixing Peng, Detao Bai, Shimin Yao, Boyuan Sun, Xiang Chen, Shenghao Fu, Xihan Wei, Liefeng Bo, et al. Humanomni: A large vision-speech language model for human-centric video understanding. *CoRR*, 2025. 3, 5, 6

[31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 3

[32] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 3

[33] Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025. 3

[34] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025.

[35] Sungjune Park, Hyunjun Kim, Junho Kim, Seongho Kim, and Yong Man Ro. Dip-r1: Deep inspection and perception with rl looking through and understanding complex scenes. *arXiv preprint arXiv:2505.23179*, 2025.

[36] Liang Chen, Lei Li, Haozhe Zhao, Yifan Song, and Vinci. R1-v: Reinforcing super generalization ability in vision-language models with less than $3. `https://github.com/Deep-Agent/R1-V`, 2025. Accessed: 2025-02-02.

[37] Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu Zeng, Yinan He, Yali Wang, Yu Qiao, Yi Wang, and Limin Wang. Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*, 2025.

[38] Hosu Lee, Junho Kim, Hyunjun Kim, and Yong Man Ro. Refocus: Reinforcement-guided frame optimization for contextual understanding. *arXiv preprint arXiv:2506.01274*, 2025.

[39] Zhenghao Xing, Xiaowei Hu, Chi-Wing Fu, Wenhai Wang, Jifeng Dai, and Pheng-Ann Heng. Echoink-r1: Exploring audio-visual reasoning in multimodal llms via reinforcement learning. *arXiv preprint arXiv:2505.04623*, 2025. 3

[40] Tong Wei, Yijun Yang, Junliang Xing, Yuanchun Shi, Zongqing Lu, and Deheng Ye. Gtr: Guided thought reinforcement prevents thought collapse in rl-based vlm agent training. *arXiv preprint arXiv:2503.08525*, 2025. 3

[41] Yi Chen, Yuying Ge, Rui Wang, Yixiao Ge, Junhao Cheng, Ying Shan, and Xihui Liu. Grpo-care: Consistency-aware reinforcement learning for multimodal reasoning. *arXiv preprint arXiv:2506.16141*, 2025. 3

[42] Qize Yang, Shimin Yao, Weixuan Chen, Shenghao Fu, Detao Bai, Jiaxing Zhao, Boyuan Sun, Bowen Yin, Xihan Wei, and Jingren Zhou. Humanomniv2: From understanding to omni-modal reasoning with context. *arXiv preprint arXiv:2506.21277*, 2025. 3

[43] Sudha Velusamy, Hariprasad Kannan, Balasubramanian Anand, Anshul Sharma, and Bilva Navathe. A method to infer emotions from facial action units. In *2011 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 2028–2031. IEEE, 2011. 3

[44] Yin Chen, Jia Li, Shiguang Shan, Meng Wang, and Richang Hong. From static to dynamic: Adapting landmark-aware image models for facial expression recognition in videos. *IEEE Transactions on Affective Computing*, 2024. 6

[45] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 5

[46] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023. 5

[47] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019. 5

[48] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 5

# Learning What to Attend First: Modality-Importance-Guided Reasoning for Reliable Multimodal Emotion Understanding

## Supplementary Material

## 7. Experimental Setup Details

### 7.1. Datasets

**EMER** [27] is a small, high-quality emotional video dataset containing 332 samples, each manually verified by humans for both reasoning and answer correctness. It is primarily used during the cold-start SFT stage of MLLM training to initialize the model's reasoning capability. This dataset contains five emotions: angry, sad, surprise, worried, and happy, and provides corresponding textual descriptions for each sample.

**DFEW** [7] is a large-scale facial expression video dataset that contains 9,362 training samples and covers seven basic emotions: angry, happy, surprise, disgust, sad, fear and neutral. The dataset additionally provides 2,342 test samples for evaluation.

**MAFW** [6] is a multimodal emotion dataset that contains 7,341 training samples. This dataset includes seven basic emotions as well as four compound categories, namely anxiety, contempt, disappointment and helplessness. It additionally provides 1,831 test samples.

### 7.2. Evaluation metrics

To assess the accuracy of emotion recognition, we adopt two commonly used metrics by following prior works [11, 14]: (i) Unweighted Average Recall (UAR), which computes the average recall across all emotion categories without considering their frequencies; and (ii) Weighted Average Recall (WAR), which weights each class by its occurrence proportion to account for class imbalance. Together, these two metrics provide a balanced evaluation of the model's performance on both frequent and infrequent emotion categories. Formally, let $C$ denote the total number of emotion classes, $N_i$ the number of samples in the $i$-th class, and $TP_i$ the number of correctly predicted samples for class $i$. WAR and UAR are defined as follows:

$$\text{WAR} = \frac{\sum_{i=1}^{C} \text{TP}_i}{\sum_{i=1}^{C} N_i} \qquad (4)$$

$$\text{UAR} = \frac{1}{C} \sum_{i=1}^{C} \frac{\text{TP}_i}{N_i} \qquad (5)$$

To evaluate the emotional coherence of generated explanations, we leverage three complementary metrics by following [16]: (i) Explanation Emotion Accuracy (EEA), which measures how accurately the emotion expressed in the explanation aligns with the ground-truth emotion label; (ii) Explanation–Prediction Consistency (EPC), which quantifies the degree of consistency between the emotion reflected in the explanation and the model's predicted emotion; and (iii) Faithful Consistency Rate (FCR), which assesses whether the explanation, the predicted emotion, and the target emotion are mutually consistent.

In our experiments, we utilize GPT-4.1-mini to infer the emotion expressed in the generated explanation. Given a total of $S$ samples, for the $i$-th sample, let $y_i$ denote the ground-truth emotion, $\hat{y}_i$ the model's predicted emotion label, and $e_i$ the emotion derived from the reasoning text. The metrics are defined as:

$$\text{EEA} = \frac{1}{S} \sum_{i=1}^{S} \mathbb{1}[e_i = y_i], \qquad (6)$$

$$\text{EPC} = \frac{1}{S} \sum_{i=1}^{S} \mathbb{1}[e_i = \hat{y}_i], \qquad (7)$$

$$\text{FCR} = \frac{1}{S} \sum_{i=1}^{S} \mathbb{1}[e_i = y_i \wedge \hat{y}_i = y_i]. \qquad (8)$$

## 8. Qualitative Analysis

### 8.1. Analysis of Attention Distribution

To investigate the source of the improvements in consistency metrics, we analyze the attention distribution of the models during the reasoning process. Specifically, we calculate the relative attention scores allocated to visual tokens, audio tokens, and generated text tokens (excluding instruction prompts) at each step of token generation. The analysis was conducted on the DFEW test set, categorized into two subsets based on the MI. To account for varying lengths of generated reasoning across samples, we normalized the generation steps to a fixed length of 100. For the model consisting of 28 layers, we specifically examine the 21st layer, with attention scores averaged across all heads. Figure 5 illustrates the modality-wise attention distribution over generated tokens for MIGR, ERV, and R1-Omni.

Our analysis reveals two key observations that distinguish MIGR from existing baselines:

**Enhanced Grounding on Multimodal Evidence** As shown in Figure 5, a prominent difference is that MIGR demonstrates significantly higher peak attention scores for both visual and audio modalities compared to the baseline models, ERV and R1-Omni. This indicates that MIGR is

**(a)** Modality-wise Attention Distribution over Generated Tokens of **MIGR**

**(b)** Modality-wise Attention Distribution over Generated Tokens of **ERV**

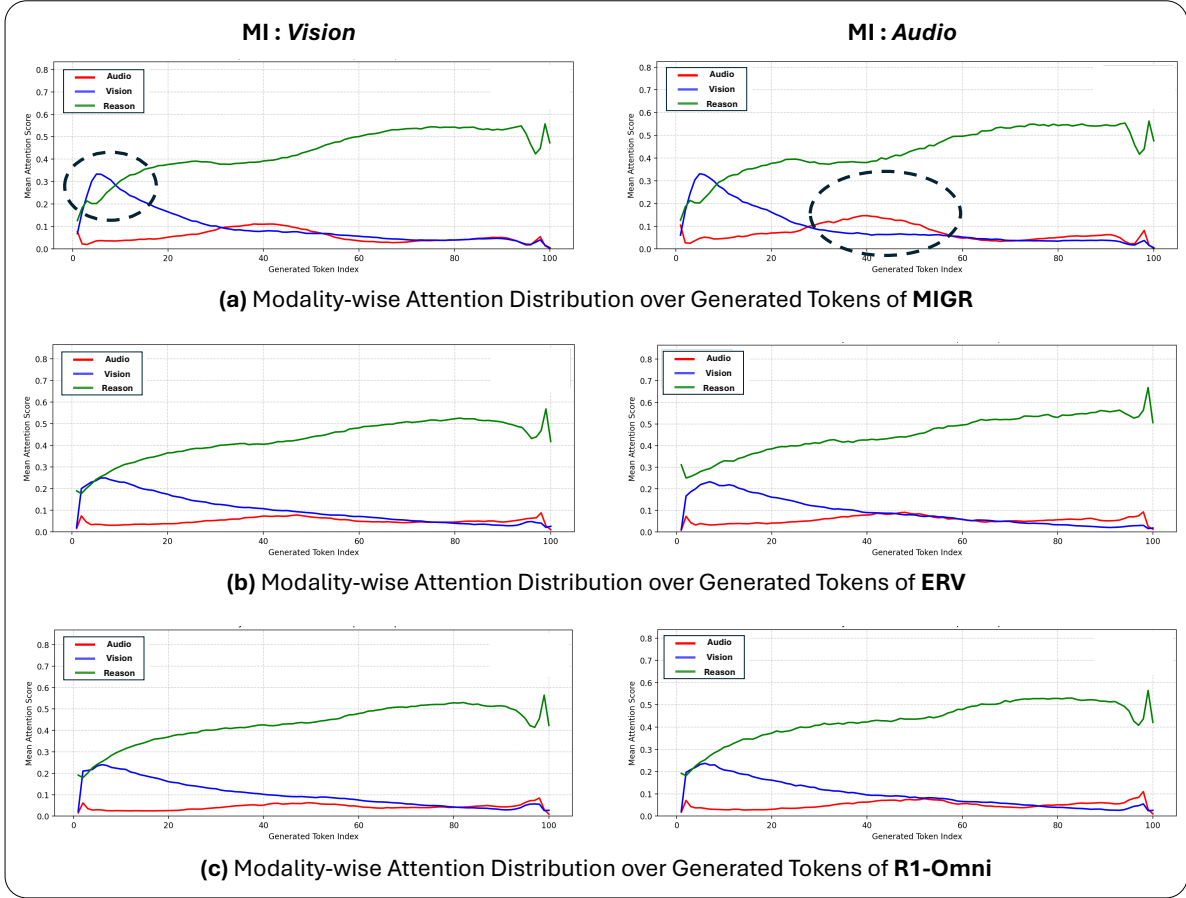**(c)** Modality-wise Attention Distribution over Generated Tokens of **R1-Omni**

Figure 5. Comparison of modality-wise attention distributions over generated tokens among MIGR, ERV, and R1-Omni.

capable of attending to multimodal evidence with greater intensity, thereby reducing reliance solely on language priors. This observation aligns with the qualitative results in Figure 6, where baseline models often fail to capture audio cues or generate hallucinations (marked in gray), whereas MIGR produces reasoning well-grounded in the input data.

**Dynamic Attention Shift based on Modality Importance** Additionally, MIGR exhibits a distinct attention pattern that adapts to the dominant modality. As highlighted by the dashed circles in Figure 5(a), when the emotion-dominant modality is Vision ($MI : Vision$), the model's attention to visual tokens is explicitly heightened. Conversely, when the Audio modality is dominant ($MI : Audio$), there is a clear surge in attention toward audio tokens. This contrasts with the baselines, which tend to exhibit a static attention bias (often towards vision or text) regardless of which modality holds the critical emotional information. This demonstrates that MIGR successfully focuses on the informative modality identified by the MI, thereby contributing to the generation of more reliable and emotion-grounded reasoning.

## 8.2. Additional Qualitative Results

Figures 6 and 7 present concrete examples that further illustrate the distinct advantages of our model compared to baseline methods.

**Precise Capture of Auditory Emotional Cues** Figure 6 demonstrates the difference in focusing capability regarding auditory cues. In case (a), despite the presence of the protagonist's clear painful speech, baseline models fail to detect it or provide reasoning with contradictory emotions, such as "excited." Conversely, in case (b), which lacks explicit emotional speech, baseline models tend to hallucinate audio descriptions corresponding to the serious visual cues. In contrast, MIGR accurately captures the "sad and tearful" audio signals in (a) and correctly infers the absence of emotion-related information in the audio of (b), demonstrating superior audio-visual discrimination.

**Alignment between Reasoning and Prediction** Figure 7 highlights the improved consistency performance of our approach. In (a), MIGR maintains the "surprise" emotion consistently from the initial reasoning phase to the final prediction. In (b), while MIGR initially considers both "happy" and "surprise," it successfully refines its judgment through
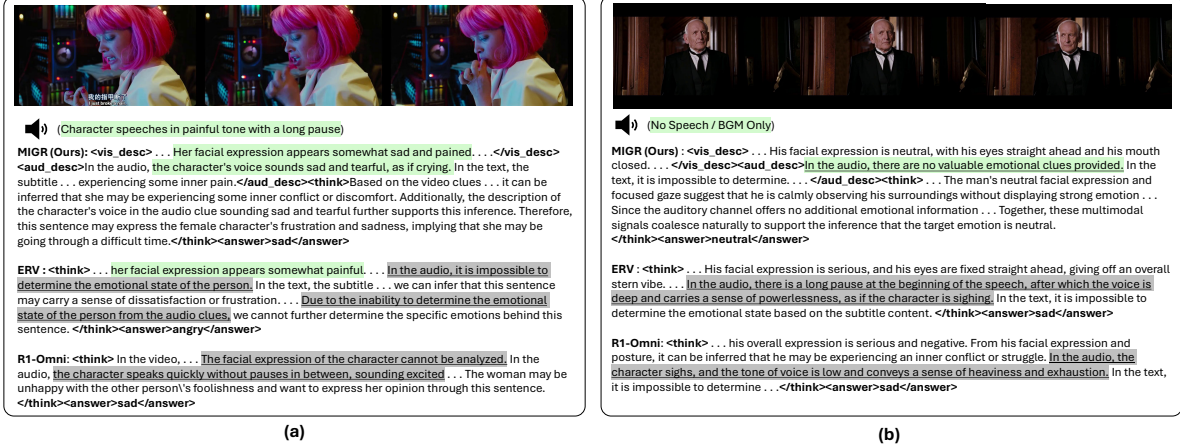
Figure 6. Qualitative results comparing MIGR, ERV, and R1-Omni. Samples (a) and (b) demonstrate MIGR's capability to precisely reason about auditory emotional cues.
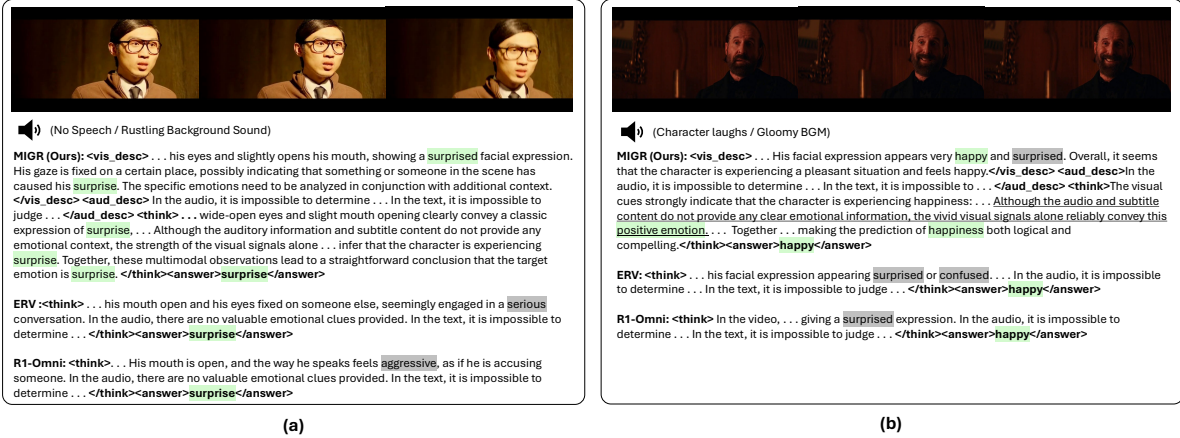


Figure 7. Qualitative results comparing MIGR, ERV, and R1-Omni. Samples (a) and (b) demonstrate MIGR's capability to consistently align its reasoning with the predicted emotional state.

a progressive reasoning process, converging to "happy" as the final prediction. On the other hand, baseline models exhibit reasoning processes that are inconsistent with their final predicted emotions. These qualitative results provide concrete examples that corroborate the consistency analysis presented in Table 2.

# 9. Additional Analysis of Performance Degradation in the Surprise Category

## 9.1. Analysis of Annotation Ambiguity

Although the DFEW dataset is utilized as a single-label benchmark, its ground truth generation involves a voting process by 10 trained human annotators. A video is assigned a single emotion label only if it receives a consensus vote of over 60%. Importantly, other emotions perceived in the video are also recorded as concurrent labels. To investigate the misclassification of Surprise samples by MIGR, we

Table 6. Analysis of **concurrent human annotations** for the 47 Surprise samples misclassified by MIGR. Although the ground truth is Surprise, a significant portion of annotators voted for **Neutral** and **Fear**, indicating the inherent ambiguity and compound nature of these samples.

| Metric | Concurrent Emotion Annotation | | | | | | Total Target |
|---|---|---|---|---|---|---|---|
| | Hap | **Neu** | Ang | Dis | **Fea** | Sad | Samples |
| **# Samples** | 0 | **35** | 9 | 4 | **23** | 5 | 47 |
| **Ratio (%)** | 0.0 | **74.5** | 19.1 | 8.5 | **48.9** | 10.6 | 100.0 |

analyzed the distribution of these concurrent annotations for the 47 samples where MIGR's prediction diverged from the ground truth. As detailed in Table 6, we found that a significant portion of these samples received annotator votes for the emotions predicted by MIGR. Specifically, 48.9% of these samples contained annotations for Fear, and 74.5% included votes for Neutral. This high overlap between the
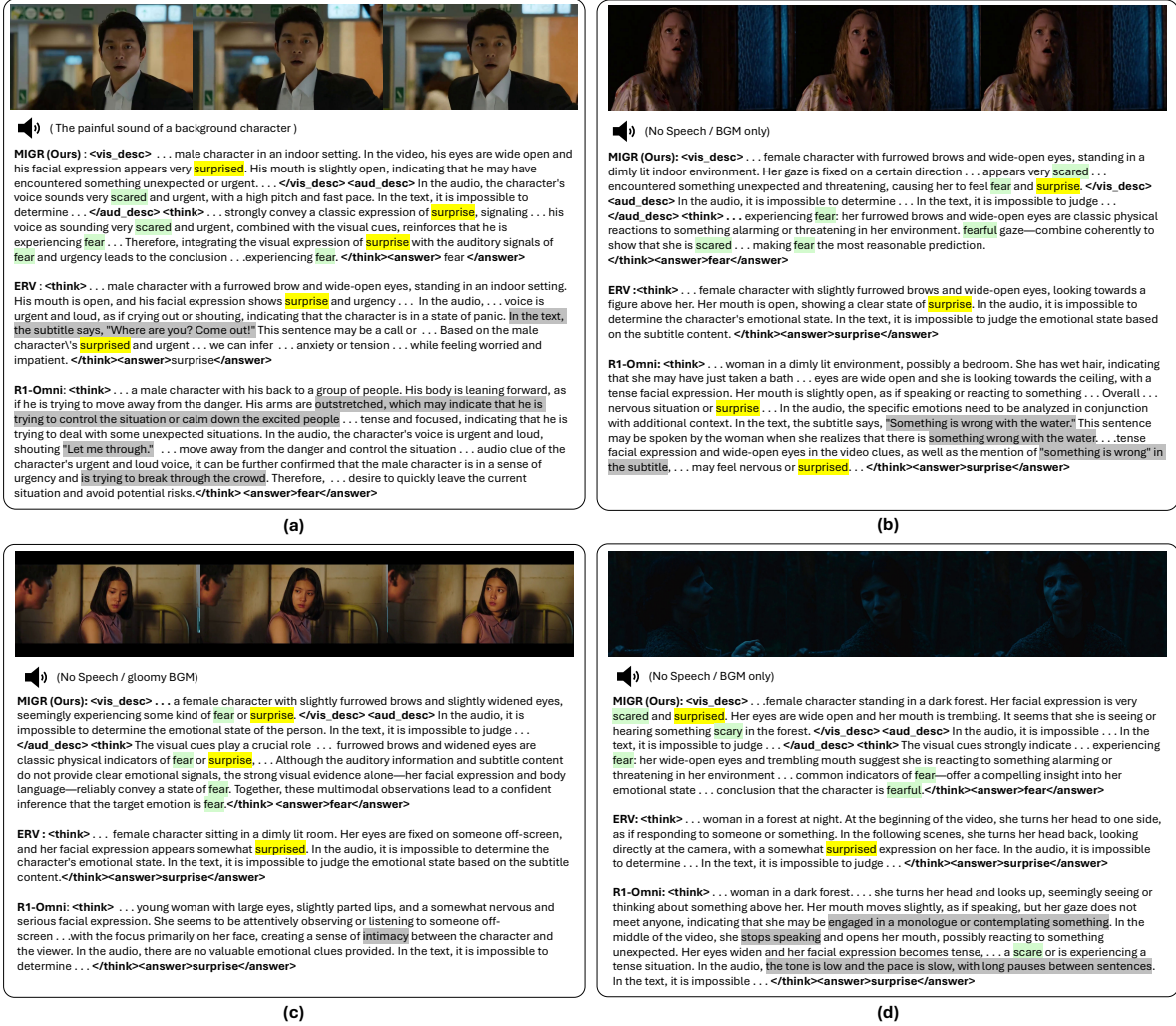
Figure 8. Qualitative comparison of samples with Surprise as the ground truth. While MIGR misclassifies the final label as Fear, its reasoning process accurately captures the compound emotional state of both Surprise and Fear. In contrast, ERV focuses solely on a single emotion in its reasoning, and R1-Omni exhibits hallucinations by basing its inference on non-existent subtitles. Text highlighted in yellow and green corresponds to reasoning associated with Surprise and Fear, respectively, while gray highlights indicate hallucinations regarding content absent from the input video.

concurrent human annotations and MIGR's predictions (Table 4) suggests that MIGR is not merely making erroneous predictions; rather, it is sensitive to the subtle, multi-faceted nature of emotions that even human annotators perceive differently.

## 9.2. Qualitative Analysis of Performance Degradation

Building on this insight, we revisited the qualitative performance discussed in Sections 5.1.3 and 5.1.4. While a notable portion of Surprise samples were classified as Fear, our analysis shows that MIGR successfully captures the Surprise element within its reasoning process in 83% of these cases (as shown in Table 4).

Consistent with the previous annotation analysis, our qualitative inspection reveals that these misclassified samples often exhibit compound emotions, where facial cues. The generated reasoning text confirms that MIGR accurately recognizes this complexity, explicitly discussing the presence of both Fear and Surprise elements.

In contrast, while the ERV often matched the ground truth label, its reasoning tended to focus exclusively on the Surprise features, failing to capture the co-occurring Fear signals or the emotional ambiguity. Furthermore, we observed critical hallucination issues in baselines. For example, as shown in Figure 8(a), in a scene containing only background sound, ERV and R1-Omni generated reasoning based on non-existent dialogue (e.g., "Where are you?

Come Out!"). This indicates a susceptibility to text-bias, leading to fabricated evidence. Conversely, MIGR effectively mitigates such hallucinations and provides robust reasoning by anchoring its output to the actual multimodal evidence, attributed to our Modality-Grounded Reasoning objective.