OPEN-VOCABULARY MULTIMODAL EMOTION RECOG NITION: DATASET, METRIC, AND BENCHMARK

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal Emotion Recognition (MER) is an important research topic. This paper advocates for a transformative paradigm in MER. The rationale behind our work is that current approaches often rely on a limited set of basic emotion labels, which do not adequately represent the rich spectrum of human emotions. These traditional and overly simplistic emotion categories fail to capture the inherent complexity and subtlety of human emotional experiences, leading to limited generalizability and practicality. Therefore, we propose a new MER paradigm called Open-vocabulary MER (OV-MER), which encompasses a broader range of emotion labels to reflect the richness of human emotions. This paradigm relaxes the label space, allowing for the prediction of arbitrary numbers and categories of emotions. To support this transition, we provide a comprehensive solution that includes a newly constructed database based on LLM and human collaborative annotations, along with corresponding metrics and a series of benchmarks. We hope this work advances emotion recognition from basic emotions to more nuanced emotions, contributing to the development of emotional AI.

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Research on emotions has a history spanning two centuries. As early as the 19th century, Charles
Darwin conducted pioneering research about the evolutionary origins and possible purposes of
emotions, explaining the emotional expressions of humans and animals (Darwin, 1872). In 1884,
William James revealed the process of emotion generation, noting that stimuli trigger activities in
the autonomic nervous system, which in turn produces an emotional experience in the brain (James, 1884). With the rapid development of AI, emotions have garnered increasing attention. Minsky
(1988) highlighted the importance of emotions: *The question is not whether intelligent machines can have any emotions, but whether machines can be intelligent without any emotions.*

Humans convey emotions through various modalities, which gives rise to the task of Multimodal
Emotion Recognition (MER) (Lian et al., 2024b). The basis of MER is how to model emotions.
Currently, emotion models can be broadly classified into two categories: dimensional models and
discrete models. The former uses multi-dimensional space to describe emotions, with the arousalvalence dimension being widely adopted; valence describes the pleasantness of a stimulus, while
arousal indicates the intensity of emotion provoked by a stimulus (Warriner et al., 2013). However,
these definitions require specialized psychological knowledge, making them abstract and difficult
for the ordinary person to understand. This can lead to discrepancies between different annotators,
posing challenges for subsequent applications.

Discrete models align more closely with human understanding of emotions. Ekman (1992) proposed the basic emotion theory, suggesting that there are six basic emotions: *anger*, *disgust*, *fear*, *happiness*, *sadness*, and *surprise*. This theory is widely used in MER, where researchers typically limit the label space to these basic emotions and use multiple annotators to select the most likely label through majority voting. We refer to this task as One-hot MER (OH-MER). Considering that emotions can be compound, researchers further propose Multi-label MER (ML-MER), allowing each sample to have multiple labels (Li et al., 2017). However, both OH-MER and ML-MER generally have limited label spaces and numbers. Nevertheless, Plutchik (2001) pointed out that humans can express approximately 34,000 different emotions. Therefore, restricting the label space and number will overlook some nuanced emotions. 054 One-hot MER Multi-label MER **Open-vocabulary MER** 056 ම ශ් <u>a a</u> Lahel One-hot Label: surp Une-brit Label: surprise Description: In the video, the screen shows a male character in an indoor setting. At the beginning of the video, his eyes are wide open and him mouth is also open, indicating a surprised field expression. In the following scenes, he looks around, sceningly explaining or marning something to the people around him. Overall, his emotions are not positive or optimistic. In the audio, the character speaks with a stutter, which usually expresses fielings of nervoaness, anxiety, or unesse. Combined with the text content, the character seems to be unhappy and any due to the prejudice of the people around him. The subtitie in the text says, "Why are you all looking at me like that? So, as long as it's a woman, does sha have to have a rationship with me?" This sentence expresses the male character's dissatisfaction and anger towards the people around him. Based on the suprised and negative field expression of the male character in the video clues, as well as the stuttering speech in the audio clues, we can infer that the male character is expressing a feeling of dissatisfaction and anger in this sentence. He may feel troubled by the prejudice of the people around him and sumhappy with this unfair treatment. Description - **OV labels**: suprise, nervous, dissatisfield 059 <mark>°°</mark> = Label 060 Number 061 One Label and Y instand 062 LLM 063 Manner 064 scription → OV labels: sur se, nervous, dissatisfied 065 (a) Task comparison (b) Label comparison

Figure 1: Comparison. (a) Task Comparison: We compare the differences among three tasks (one-hot MER, multi-label MER, and open-vocabulary MER) across three aspects (label space, label number, and annotation manner); (b) Label Comparison: We provide an example to visualize the one-hot and OV labels. Since the original video contains real people, we use DemoAI to remove personal information to address copyright concerns. This paper uses emotion-related descriptions as a bridge to extract OV labels. We observe that OV labels contain richer emotions.

073 074 075

066 067

In this paper, we extend traditional MER to Open-vocabulary MER (OV-MER), where we allow the 076 prediction of emotions across any number and categories. Figure 1(a) provides comparisons between 077 different tasks. To facilitate further research, we build an initial dataset, define evaluation metrics, and develop solutions: (1) For the dataset, we propose a human-LLM collaboration strategy. Com-079 pared to human-only annotation, our strategy can leverage LLM to enhance the label richness (see Section 2); (2) For the metrics, since there is no fixed label space, the model may predict closely 081 related but differently expressed emotions. For example, the ground truth is *joyful* and the model 082 predicts happy. To provide more reliable evaluation results, we first group similar emotions and 083 specifically design metrics for this task (see Section 3); (3) For the solutions, traditional classifiers 084 rely on fixed label spaces. However, OV-MER does not restrict the label space or number of labels, 085 necessitating the definition of new solutions (see Section 4).

A natural question arises: why is open vocabulary so important for MER? The reason is that OV can generate more nuanced emotions, leading to more accurate and reliable MER. As illustrated in Figure 1(b), labeling an emotion solely as *surprise* is not sufficiently informative—we cannot ascertain whether the *surprise* is positive or negative, which limits its practicality, such as in humancomputer interaction (HCI) applications. In contrast, our OV-MER provides emotions like *surprise* along with additional descriptors such as *nervous* and *dissatisfied*, offering a more comprehensive and insightful understanding of the emotional state. Therefore, OV-MER facilitates the transition from basic to nuanced emotion recognition, advancing the development of emotion AI. Appendix A provides more detailed motivation. In summary, we make the following key contributions:

095

098

099

102

- 096 097
- **Paradigm**. We propose a new paradigm in MER called OV-MER. This paradigm transitions from traditional MER to a framework that enables the prediction of any number and category of emotions, thereby advancing emotion AI toward real-world applicability by capturing the full spectrum of human emotions.
- **Groundwork**. We lay the groundwork for OV-MER by constructing datasets, defining evaluation metrics, and proposing effective solutions. Our dataset enhances label richness through human-LLM collaboration. Meanwhile, we introduce new evaluation metrics that leverage emotional relevance to achieve more reliable results.
- Benchmark. We establish zero-shot benchmarks for OV-MER by conducting extensive experiments and providing detailed analysis. This task can serve as an important evaluation benchmark for multimodal LLMs (MLLMs), challenging their ability to integrate multimodal clues and capture subtle temporal variations in emotional expression.



129

130

131

132

133

134 135 136

137 138

139

140

141

142

143

144

145

Figure 2: Dataset construction. (a) CLUE-Multi Generation: For audio and video, we use ALLM and VLLM to extract initial clues, followed by two rounds of manual checks to eliminate errors and duplicates while adding missing content. Each round involves multiple annotators, with no overlap between annotators in the two rounds. Finally, we merge the checked clues with text to generate CLUE-Multi. (b) Ground-truth OV Label Extraction: There are certain differences in the labels extracted from different languages. To eliminate language influence and achieve consensus labels, we merge these labels and conduct manual checks.

2 THE OV-MERD DATASET CONSTRUCTION

Although the concept of OV-MER is intuitive and holds great promise, its practical implementation faces significant challenges. The main difficulty lies in the broad and subtle range of human emotions, making comprehensive labeling a complex task. Traditional annotation methods are limited by their predefined emotion categories, which are often insufficient for the needs of OV-MER. In Figure 2, we propose a human-LLM collaboration strategy that consists of two steps: CLUE-Multi generation and emotion label extraction. Ultimately, we create a dataset called OV-MERD, which offers a richer set of emotions compared to existing datasets (see Table 1). This dataset is an extension of MER2023 (Lian et al., 2023). Specifically, MER2023 is collected from movies and TV series, and we randomly select a portion of data for further annotation to build our OV-MERD. 146

147 148

149

2.1 CLUE-MULTI GENERATION

150 During the annotation process, we observe that descriptions generated through human-LLM collab-151 oration are more detailed. This indicates that when annotations are solely performed by humans, 152 they tend to focus on major clues while neglecting minor ones (see Section 5 for more analysis). In 153 this section, we provide a detailed overview of our human-LLM collaboration strategy.

154

155 **Pre-annotation.** Initially, we attempt to annotate visual and acoustic clues directly. However, the 156 descriptions obtained in this way cannot cover all information. Therefore, we explore using other 157 models for pre-annotation. (1) For video, given the strong visual understanding capabilities of GPT-158 4V ("gpt-4-vision-preview"), we use it as video LLM (VLLM) for pre-annotation. Since GPT-4V 159 only supports image input, we uniformly sample three frames from each video and input them into GPT-4V. We discuss the reasons for sampling three frames in Appendix I. (2) For audio, we use the 160 open-source SALMONN (Tang et al., 2023) as audio LLM (ALLM) for pre-annotation, as GPT-4V 161 does not support audio input.

163	Table 1: Dataset comparison. See Appendix G for a comprehensive comparison.				
164	Dataset	Modality	Annotation Type	# Categories	# Labels per Sample
10-1	MOUD (Pérez-Rosas et al., 2013)	A,V,T	Dimensional Emotion	1	1
165	CMU-MOSI (Zadeh et al., 2017)	A,V,T	Dimensional Emotion	1	1
166	CH-SIMS (Yu et al., 2020)	A,V,T	Dimensional Emotion	1	1
100	CH-SIMS v2 (Liu et al., 2022a)	A,V,T	Dimensional Emotion	1	1
167	SEMAINE (McKeown et al., 2011)	A,V,T	Dimensional Emotion	5	1
169	MSP-IMPROV (Busso et al., 2016)	A,V,T	Discrete Emotion	4	1
100	IEMOCAP (Busso et al., 2008)	A,V,T	Discrete Emotion	10	1
169	MELD (Poria et al., 2019)	A,V,T	Discrete Emotion	7	1
170	MER2023 (Lian et al., 2023)	A,V,T	Discrete Emotion	6	1
170	MER2024 (Lian et al., 2024a)	A,V,T	Discrete Emotion	6	1
171	OV MERD (Ours)	AVT	Discrete Emotion	248	1~9, most 2~4
172	OV-MERD (Ours)	A, V, I	Discrete Emotion	(arbitrary label)	(arbitrary number)

162

174 Manual Check. As part of our quality assurance procedures, we perform a detailed examination 175 of the pre-annotated results. For visual clues, GPT-4V may generate hallucinated responses, i.e., 176 clues that do not actually exist. Additionally, there are repeated expressions and some temporal 177 association clues are missing. Therefore, we hire annotators to eliminate errors and duplicates, as 178 well as add missing content. For acoustic clues, ALLM struggles to capture emotion-related par-179 alinguistic features. The main reason is that current ALLM primarily focuses on tasks like ASR or 180 audio event detection (Tang et al., 2023), with less emphasis on paralinguistic information. Hence, we hire multiple annotators to focus on the speaker's intonation and other emotion-related paralin-181 guistic clues. To reduce subjective bias, we conduct two rounds of manual checks, each involving 182 different annotators. These annotators are experts in affective computing and are familiar with the 183 definitions of emotions. Ultimately, these checked clues can accurately reflect the video content. Appendix J provides the annotation guideline and layout of the annotation platform. 185

CLUE-Multi Generation. Subsequently, we leverage the reasoning capabilities of LLM to merge 187 all clues. Specifically, we use GPT-3.5 ("gpt-3.5-turbo-16k-0613") as the LLM and ask it to merge 188 textual, acoustic, and visual clues to infer emotional states. The output of this process is an emotion-189 related description, denoted as *CLUE-Multi* (see Figure 2). In Appendix E, we further discuss the 190 details of this merging process and the reasons behind it. Overall, the above annotation pipeline 191 reflects the collaboration between humans and LLMs. In Section 5, we compare human-only anno-192 tation with this strategy, observing that the latter can generate richer descriptions.

193 194

195

2.2 GROUND-TRUTH OV LABEL EXTRACTION

196 Label Extraction. After that, we use the LLM to extract emotion labels from *CLUE-Multi*. This process relies on GPT-3.5, which we request to identify emotional states based on the provided 197 descriptions without restricting the label space. See Appendix E for more details. 198

199

Language Impact. We further explore the language impact. In Figure 2, we first extract OV 200 labels from English and Chinese descriptions, obtaining Y_{EE} and Y_{CC} . Then, we translate them 201 into the other language, yielding Y_{EC} and Y_{CE} . After that, we measure the similarity between 202 different sets and report results in Figure 2. In Appendix K, we detail our experimental design 203 and similarity metric. We observe that the labels extracted from different languages exhibit some 204 differences. For example, the similarity score between Y_{EE} and Y_{CE} is 0.82, which may be due 205 to the varying definitions of emotions in different languages. To eliminate language influence and 206 achieve consensus labels, we merge the labels extracted from both languages and conduct manual 207 checks. These checked labels are regarded as the ground truth.

208 209

2.3 OV-MERD DATASET

210

211 Finally, we construct a dataset called OV-MERD. This dataset is an extension of MER2023 (Lian 212 et al., 2023), from which we randomly select a portion of samples for further annotation. Table 1 213 compares OV-MERD with existing datasets. We observe that our OV-MERD dataset contains 248 emotion categories and most samples have 2 to 4 labels, far exceeding those in current datasets. In 214 Appendix F, we observe that OV-MERD encompasses a broader range of emotions, including some 215 that have been rarely discussed in previous research, such as shy, nervous, and grateful.

216 3 **EVALUATION METRIC**

217 218

221

224 225

229

230

231

232

233

234

235

236

Defining evaluation metrics for OV-MER presents significant challenges: (1) **OV-MER supports** 219 **predicting emotions of any category.** Thus, the model may predict closely related but differently 220 expressed emotions. To provide more reliable evaluation results, we first group the emotions based on their similarities. (2) OV-MER allows for the prediction of an arbitrary number of labels. 222 Thus, traditional evaluation metrics designed for a fixed number of labels may not be applicable. In this section, we propose set-based evaluation metrics specifically tailored for this task. 223

3.1 GROUPING

226 We propose two grouping strategies: one based on GPT and the other based on the emotion wheel 227 (EW) Plutchik (1980). In the following experiments, we default to using GPT-based grouping. 228

GPT-based Grouping. The most direct approach is to use GPT-3.5 to group all labels based on their similarity: Please assume the role of an expert in the field of emotions. We provide a set of emotions. Please group the emotions, with each group containing emotions with the same meaning. Directly output the results. The output format should be a list containing multiple lists. However, the evaluation results may be affected by the API version. For example, if OpenAI deprecates an old API, the results based on that API will become difficult to reproduce. Additionally, this process is costly (see Appendix P). Therefore, we attempt to find a replacement for GPT-based grouping.

237 EW-based Grouping. EW is a psychological model that catego-238 rizes emotions in a structured manner. The inner part shows core 239 emotions while moving to the outer part reveals more nuanced emotions. Therefore, EW naturally provides emotion grouping informa-240 tion. Since there is no consensus on EW, we select five representa-241 tive wheels $W1 \sim W5$. See Appendix R for details. 242

243 Take W1 as an example (see Figure 3). Before calculating the met-244 rics, we define some symbols. We group the labels by their lev-245 els from the innermost to the outermost as $L_{w_1}^1$, $L_{w_1}^2$, and $L_{w_1}^3$. Specifically, $L_{w_1}^1 = \{ mad, \cdots, scared \}$. Next, we define a func-246 tion $m_{w_1}^{i \to j}(\cdot)$ that maps the labels in $L_{w_1}^i$ to the corresponding labels in $L_{w_1}^j$. From inner to outer (i < j), $m_{w_1}^{i \to j}(\cdot)$ is a many-to-one mapping; from outer to inner (i > j), $m_{w_1}^{i \to j}(\cdot)$ is a one-247 248 249 250 to-many mapping. For example, $m_{w_1}^{2 \to 1}(\text{angry}) = \{\text{mad}\}$ and



Figure 3: Emotion wheel W1.

251 $m_{w_1}^{1\to 2}(\text{mad}) = \{\text{offended}, \cdots, \text{humiliated}\}$. We collect all the labels from these emotion wheels 252 and represent them as EW, i.e., $\{L_{w_i}^j, 1 \le i \le 5, 1 \le j \le 3\}$. We denote the labels in EW as y_w . 253

Considering that the emotional categories in EW are still limited, we perform some label expansion 254 operations. Specifically, we repeatedly call GPT-3.5, asking it to generate synonyms for each label. 255 The prompt used is as follows: Please retrieve the synonyms for the following words and output 256 them in a table format. Then, we generate EW-S, i.e., $\{f(y_w) = \{y_f^1, ..., y_f^n\}, y_w \in EW\}$, where 257 $f(\cdot)$ is a function that maps each label y_w to its synonym y_f . We also define its inverse function 258 $f'(\cdot)$, which maps different synonyms y_f back to their base label y_w . 259

To eliminate the influence of word forms (e.g., happy and happiness), we further ask GPT-3.5 260 multiple times to generate different forms for each label. The prompt used is as follows: Please 261 output different forms of the following word in a list format. After that, we obtain EW-SF, i.e., 262 $\{g(y_f) = \{y_q^1, ..., y_q^m\}, y_f \in \text{EW-S}\}, \text{ where } g(\cdot) \text{ is a function that maps each label } y_f \text{ to its differ-}$ 263 ent forms y_g . We also define its inverse function $g'(\cdot)$, which maps different labels y_g back to their 264 base form y_f . Finally, we define different types of metrics: 265

- (1) M1. We use $g'(\cdot)$ to map all labels to their corresponding y_f . 266
- 267 (2) M2. We use $f'(g'(\cdot))$ to map all labels to their corresponding y_w . 268
- (3) M3. We use the emotion wheel during metric calculation. Specifically, we first use $f'(q'(\cdot))$ to 269 map all labels to their corresponding y_w . Then, we define two grouping functions, L1 and L2. For

L1, we map all labels to their corresponding $L_{w_i}^1$:

$$\begin{cases} y_w, & \text{if } y_w \in L^1_{w_i} \\ m^{2 \to 1}_{w_i}(y_w), & \text{if } y_w \in L^2_{w_i} \\ m^{2 \to 1}_{w_i}(m^{3 \to 2}_{w_i}(y_w)), & \text{if } y_w \in L^3_{w_i} \end{cases}$$
(1)

For L2, we map all labels to their corresponding $L_{w_i}^2$:

$$\begin{cases} \text{select one label in } m_{w_i}^{1 \to 2}(y_w), & \text{if } y_w \in L^1_{w_i} \\ y_w, & \text{if } y_w \in L^2_{w_i} \\ m_{w_i}^{3 \to 2}(y_w), & \text{if } y_w \in L^3_{w_i} \end{cases}$$
(2)

3.2 METRIC DEFINITION

Then, we convert the above emotion grouping information into a function $G(\cdot)$, which can map each label to its group ID. Specifically, suppose $\{y_i\}_{i=1}^M$ and $\{\hat{y}_i\}_{i=1}^N$ are the ground truth and predictions, where M and N are the number of labels. We first map each label into its group ID:

$$\mathcal{Y} = \{G(x) | x \in \{y_i\}_{i=1}^M\}, \hat{\mathcal{Y}} = \{G(x) | x \in \{\hat{y}_i\}_{i=1}^N\}.$$
(3)

Finally, we design set-based metrics for performance evaluation. Specifically, Precision_s indicates the number of correctly predicted labels; Recall_s indicates whether the prediction covers all ground truth; and F_s is the harmonic mean of two metrics, which is used for the final ranking. It is important to note that changing the label order in \mathcal{Y} and $\hat{\mathcal{Y}}$ does not result in any change in performance.

$$\operatorname{Precision}_{S} = \frac{|\mathcal{Y} \cap \hat{\mathcal{Y}}|}{|\hat{\mathcal{Y}}|}, \ \operatorname{Recall}_{S} = \frac{|\mathcal{Y} \cap \hat{\mathcal{Y}}|}{|\mathcal{Y}|}, \ \operatorname{F}_{S} = 2 \times \frac{\operatorname{Precision}_{S} \times \operatorname{Recall}_{S}}{\operatorname{Precision}_{S} + \operatorname{Recall}_{S}}.$$
(4)

4 BASELINES FOR OV-MER

4.1 CLUE GENERATION

Figure 2 illustrates the generation process of CLUE-Multi, where we combine text with manually checked visual and acoustic clues. In this section, we further introduce the following variants.

CLUE-A/T/V. To reveal the modality impact, we propose three variants of CLUE-Multi: *CLUE-Audio*, *CLUE-Text*, and *CLUE-Video*. Their generation process is illustrated in Figure 4. (1) *CLUE-Audio*: We observe that ALLM cannot fully leverage the text, and using an additional LLM to emphasize the text can further improve performance, which is also verified in Section 5. Therefore, we merge the checked acoustic clues with text using an additional LLM; (2) *CLUE-Text*: We only use the text to infer emotional states; (3) *CLUE-Video*: Since the visual content does not contain audio and text, we only use the checked visual clues. See Appendix M for more examples.

CLUE-MLLM. MLLMs can address various multimodal tasks. Since emotion recognition relies
 on temporal information, we choose models that support at least video or audio. To generate *CLUE-MLLM*, we first use ALLM or VLLM to extract emotion-related descriptions, and then combine
 these descriptions with text using LLM. Compared with CLUE-Multi, this process does not use
 manually checked clues. Appendix N provides model cards and relevant prompts. This paper aims to
 build a zero-shot benchmark for OV-MER, without the training process. All models are implemented
 in PyTorch, and all inference processes are carried out using a 32G NVIDIA Tesla V100 GPU.

316317 4.2 METRIC CALCULATION

As shown in Figure 2, there are certain differences in the labels extracted from different languages.
Therefore, we report the results for both English and Chinese descriptions. In Figure 4, for the
Chinese branch, we first extract OV labels and then translate them into English; for the English
branch, we directly extract OV labels. Finally, we compute the evaluation metrics with the ground
truth. It is worth noting that the OV labels extracted from the monolingual CLUE-Multi differ
from the ground truth. Our ground truth combines the labels extracted from different languages and
undergoes further manual checks (see Figure 2).



Figure 4: Baselines. (a) **Preliminary:** We begin by defining some preliminary symbols. (b) **CLUE Generation:** CLUE-Video and CLUE-Audio use manually-checked clues; CLUE-Text relies solely on text; CLUE-MLLM does not involve manual checks and directly uses the outputs from ALLM or VLLM. (c) **Metric Calculation:** We rely on CLUE to predict emotion labels. Due to variations in labels extracted from different languages, we report results across different languages.

Table 2: Baseline results.	Figure 4 shows t	he metric cal	culation process

M. 1.1	T	17			English			Chinese	
Model	L	v	A	$F_{S}\uparrow$	Precision _s ↑	$\text{Recall}_{s} \uparrow$	$F_S \uparrow$	Precision _s ↑	$\text{Recall}_{S} \uparrow$
Heuristic Baseline									
Random	×	×	×	$17.42_{\pm 0.01}$	$24.85_{\pm 0.15}$	$13.42_{\pm 0.04}$	$16.59_{\pm 0.00}$	$24.70_{\pm 0.00}$	$12.48_{\pm 0.00}$
					CLUE-ML	.LM			
Qwen-Audio		×		$38.13_{\pm 0.05}$	$49.42_{\pm 0.18}$	$31.04_{\pm 0.00}$	$41.14_{\pm 0.07}$	$53.71_{\pm 0.00}$	$33.34_{\pm 0.09}$
OneLLM		×		$42.84_{\pm 0.06}$	$45.92_{\pm 0.05}$	$40.15_{\pm 0.06}$	$46.17_{\pm 0.02}$	$52.07_{\pm 0.06}$	$41.47_{\pm 0.08}$
Otter			\times	$43.51_{\pm 0.09}$	$50.71_{\pm 0.10}$	$38.09_{\pm 0.09}$	$46.22_{\pm 0.01}$	$52.65_{\pm 0.16}$	$41.18_{\pm 0.08}$
Video-LLaMA			\times	$44.73_{\pm 0.14}$	$44.14_{\pm 0.13}$	$45.34_{\pm 0.15}$	$47.26_{\pm 0.03}$	$47.98_{\pm 0.07}$	$46.56_{\pm 0.01}$
VideoChat			×	$45.53_{\pm 0.11}$	$42.90_{\pm 0.27}$	$48.49_{\pm 0.10}$	$45.57_{\pm 0.03}$	$47.20_{\pm 0.12}$	$44.05_{\pm 0.05}$
SECap		×		$45.72_{\pm 0.09}$	$54.52_{\pm 0.15}$	$39.37_{\pm 0.05}$	$45.57_{\pm 0.13}$	$55.55_{\pm 0.23}$	$38.64_{\pm 0.08}$
PandaGPT				$45.89_{\pm 0.20}$	$50.03_{\pm 0.01}$	$42.38_{\pm 0.33}$	$47.33_{\pm 0.04}$	$53.01_{\pm 0.08}$	$42.75_{\pm 0.11}$
Video-LLaVA			\times	$47.07_{\pm 0.16}$	$48.58_{\pm 0.02}$	$45.66_{\pm 0.29}$	$49.21_{\pm 0.06}$	$53.95_{\pm 0.03}$	$45.23_{\pm 0.13}$
SALMONN		×		$47.96_{\pm 0.04}$	$50.20_{\pm 0.04}$	$45.92_{\pm 0.04}$	$48.24_{\pm 0.03}$	$52.24_{\pm 0.00}$	$44.82_{\pm 0.05}$
VideoChat2			×	$49.07_{\pm 0.26}$	$54.72_{\pm 0.41}$	$44.47_{\pm 0.15}$	$48.86_{\pm 0.05}$	$57.12_{\pm 0.08}$	$42.68_{\pm 0.04}$
Video-ChatGPT			\times	$50.52_{\pm 0.06}$	$54.03_{\pm 0.04}$	$47.44_{\pm 0.07}$	$54.73_{\pm 0.00}$	$61.15_{\pm 0.10}$	$49.52_{\pm 0.06}$
OneLLM			\times	$50.52_{\pm 0.07}$	55.93 _{±0.09}	$46.06_{\pm 0.06}$	$51.44_{\pm 0.08}$	$56.43_{\pm 0.04}$	$47.26_{\pm 0.11}$
LLaMA-VID			×	$51.25_{\pm 0.09}$	$52.71_{\pm 0.18}$	$49.87_{\pm 0.00}$	$52.01_{\pm 0.02}$	$57.30_{\pm 0.00}$	$47.61_{\pm 0.03}$
mPLUG-Owl			×	$52.73_{\pm 0.13}$	$54.54_{\pm 0.13}$	$51.04_{\pm 0.13}$	$50.95_{\pm 0.06}$	$56.40_{\pm 0.11}$	$46.47_{\pm 0.18}$
Chat-UniVi			×	$53.08_{\pm 0.01}$	$53.68_{\pm 0.00}$	$52.50_{\pm 0.02}$	$53.86_{\pm 0.02}$	$58.54_{\pm 0.01}$	$49.86_{\pm 0.03}$
GPT-4V			\times	55.51 _{±0.05}	$48.52_{\pm 0.07}$	$64.86_{\pm 0.00}$	57.21 ±0.01	$54.61_{\pm 0.02}$	$60.07_{\pm 0.01}$
				_	CLUE-M/A	J/T/V			
CLUE-Text		\times	\times	$46.00_{\pm 0.06}$	$54.41_{\pm 0.15}$	$39.84_{\pm 0.01}$	$43.11_{\pm 0.25}$	$50.69_{\pm 0.26}$	$37.50_{\pm 0.23}$
CLUE-Video	×		\times	$60.55_{\pm 0.13}$	$63.29_{\pm 0.08}$	$58.05_{\pm 0.16}$	$61.73_{\pm 0.10}$	$66.47_{\pm 0.13}$	$57.62_{\pm 0.08}$
CLUE-Audio		×		$65.35_{\pm 0.04}$	$67.54_{\pm 0.08}$	$63.30_{\pm 0.00}$	$68.56_{\pm 0.07}$	$70.10_{\pm 0.06}$	$67.07_{\pm 0.08}$
CLUE-Multi				80.05 _{±0.24}	$80.03_{\pm 0.37}$	$80.07_{\pm 0.10}$	85.16 _{±0.03}	87.09 ±0.00	$83.31_{\pm 0.05}$

5 RESULTS AND DISCUSSION

In this section, we default to using GPT-based grouping and employ GPT-3.5 ("gpt-3.5-turbo-16k-0613") as LLM. We generally report evaluation results in both languages, but if no specific language is mentioned, we default to reporting results for the English branch. To mitigate the impact of randomness, we conduct each experiment twice and report the average scores and standard deviations.

Main Results on CLUE-M/A/T/V. For CLUE-M/A/T/V, most baselines use manually checked
 clues, which serve as performance upper bounds of different modality combinations. In Table 2, we
 observe that CLUE-Multi performs the best, highlighting the importance of multimodal information
 in emotion recognition. In contrast, CLUE-Text performs the worst. This is because our OV-MERD
 dataset is derived from MER2023, where the contribution of text is smaller compared to audio and
 video (Lian et al., 2024b). Relying solely on text makes it difficult to recognize emotions accurately.



Table 3: Performance of combinations of different MLLMs

Figure 5: Human-only annotation (H) vs. Human-LLM (H+L) collaboration.

410 Main Results on CLUE-MLLM. Table 2 presents the results of CLUE-MLLM. Additionally, we 411 introduce a heuristic baseline called Random, where we randomly select a label from basic emo-412 tions. This baseline reflects the lower bound of performance. We observe that MLLM generally outperforms Random, indicating that MLLM can partially address OV-MER. However, the perfor-413 mance of MLLM remains unsatisfactory, highlighting the limitations of existing MLLMs and the 414 challenges of OV-MER. Furthermore, models that perform well in Chinese often perform well in 415 English. These results suggest that the impact of language differences on rankings is limited. 416

417 Effectiveness of Multimodal Fusion. In Table 3, we select the best-performing ALLMs and 418 VLLMs and explore whether their combinations can lead to better performance. To fuse differ-419 ent modalities, we input prelabeled acoustic and visual clues into LLM and leverage its reasoning 420 capabilities for multimodal integration (see Figure 4). This approach is consistent with the fusion 421 method used during our dataset construction process (see Section 2). We observe that multimodal 422 results generally outperform ALLM-only or VLLM-only results. The reason lies in that emotions 423 are conveyed through various modalities. By integrating different modalities, we can obtain a more 424 comprehensive understanding of emotions, leading to better performance in OV-MER.

425

407

408 409

378

426 Human-only vs. Human-LLM Collaboration. To verify the effectiveness of our human-LLM 427 collaborative strategy, we additionally introduce a baseline using human-only annotation. In Figure 428 5, we compare two strategies from three aspects: the length distribution of generated clues, the distribution of label counts, and the word cloud. We observe that through human-LLM collaboration, 429 we can obtain longer descriptions, generate a broader range of emotions, and provide more diverse 430 labels for each sample. These results show that human-only annotation generally focuses on primary 431 emotions while neglecting minor ones. With the pre-annotation and semantic reasoning capabilities

447 448 449

450

451

452 453

466



Figure 6: Performance comparison of different strategies for generating CLUE-MLLM.

of LLMs, we can obtain richer emotional labels. These results validate the effectiveness of our human-LLM collaborative strategy. Meanwhile, these results suggest that the LLM-driven approach does not lead to a narrow or biased interpretation of emotions, but rather helps uncover more subtle emotional nuances. Additional experiments are provided in the Appendix O.

454 Ablation Study on CLUE-MLLM. In this section, we reveal the im-455 pact of different CLUE-MLLM generation strategies. Figure 7 introduces three methods: 1) S0 does not use text and inputs the video into 456 MLLM; 2) S1 inputs both text and video into MLLM; 3) S2 first uses 457 MLLM to extract descriptions and then combines with text using an-458 other LLM, same with the strategy in Figure 4. In Figure 6, S1 and 459 S2 generally outperform S0, indicating the importance of the text con-460 tent in OV-MER. Moreover, S2 typically performs better than S1. The 461 reason is that inputting video and text into the MLLM simultaneously 462 increases the task difficulty, and current MLLMs may struggle to han-463 dle complex prompts. S2 divides this process into two steps, reducing 464 task complexity and achieving better performance. Therefore, we adopt 465 S2 as the default strategy. More results are provided in Appendix N.

467 GPT-based vs. Matching-based Metrics. In 468 Table 2, CLUE-Multi performs the best. This 469 leads to a hypothesis: Do sentences that are more 470 "similar" to CLUE-Multi yield better emotion 471 recognition performance? The most common way to measure "similarity" is through matching-472 based metrics, with BLEU₁, BLEU₄, METEOR, 473 and ROUGE₁ being the most widely used. There-474 fore, we use CLUE-MLLM as input and calcu-475 late both GPT-based and matching-based metrics 476 (see Figure 8(a)), followed by calculating their 477 PCC scores (see Figure 8(b)). From the experi-478 mental results, we have some interesting obser-479 vations. First, the same metric across differ-480 ent languages typically exhibits high correlations. 481 However, the correlation between GPT-based and



Figure 7: Ablation study.



(a) Metric Calculation (b) Correlation Analysis

Figure 8: GPT- vs. Matching-based metrics.

matching-based metrics is not strong. For example, the highest PCC score between "F(E)" and
matching-based metrics is 0.77. These results highlight the limitations of using matching-based
metrics to evaluate the OV-MER task. The main reason is that matching-based metrics focus on
low-level word-level matches, while emotion understanding is a relatively complex and high-level
perceptual task. In Appendix Q, we provide more examples for clarification.

487	Table 4: GPT-based	IVS. EW-D	asea group	ing. we ca	iculate the	PCC score	to reveal t	neir correlat
488	Model	CPT	M1	M2	M3-	-W1	M3-	-W2
-100	Widdel	OFI	1011	IVIZ	L1	L2	L1	L2
489	Qwen-Audio	$38.13_{\pm 0.05}$	$20.49_{\pm 0.01}$	$23.37_{\pm 0.01}$	$43.85_{\pm 0.03}$	$26.60_{\pm 0.01}$	$41.52_{\pm 0.28}$	$26.68_{\pm 0.01}$
400	Otter	$43.51_{\pm 0.09}$	$22.21_{\pm 0.06}$	$27.99_{\pm 0.02}$	$49.75_{\pm 0.11}$	$33.50_{\pm 0.06}$	$49.93_{\pm 0.11}$	$33.04_{\pm 0.06}$
490	Video-LLaMA	$44.73_{\pm 0.14}$	$23.56_{\pm 0.08}$	$28.39_{\pm 0.17}$	$52.90_{\pm 0.12}$	$36.08_{\pm 0.13}$	$53.60_{\pm 0.04}$	$35.33_{\pm 0.08}$
491	VideoChat	$45.53_{\pm 0.11}$	$22.15_{\pm 0.00}$	$26.24_{\pm 0.08}$	$47.79_{\pm 0.07}$	$32.64_{\pm 0.07}$	$47.76_{\pm 0.11}$	$32.14_{\pm 0.04}$
100	SECap	$45.72_{\pm 0.09}$	$26.52_{\pm 0.01}$	$32.88_{\pm 0.03}$	$52.26_{\pm 0.03}$	$37.55_{\pm 0.03}$	$52.11_{\pm 0.03}$	$37.71_{\pm 0.03}$
492	Video-LLaVA	$47.07_{\pm 0.16}$	$25.47_{\pm 0.12}$	$30.73_{\pm 0.11}$	$54.65_{\pm 0.10}$	$37.65_{\pm 0.24}$	$54.54_{\pm 0.02}$	$38.25_{\pm 0.22}$
493	SALMONN	$47.96_{\pm 0.04}$	$23.57_{\pm 0.02}$	$28.83_{\pm 0.03}$	$54.90_{\pm 0.15}$	$38.93_{\pm 0.15}$	$54.29_{\pm 0.06}$	$37.79_{\pm 0.07}$
-100	VideoChat2	$49.07_{\pm 0.26}$	$26.92_{\pm 0.09}$	$31.40_{\pm 0.10}$	$52.38_{\pm 0.13}$	$36.44_{\pm 0.11}$	$53.56_{\pm 0.13}$	$36.91_{\pm 0.11}$
494	Video-ChatGPT	$50.52_{\pm 0.06}$	$28.99_{\pm 0.04}$	$34.05_{\pm 0.05}$	$57.66_{\pm 0.04}$	$41.48_{\pm 0.09}$	$57.37_{\pm 0.00}$	$40.95_{\pm 0.08}$
105	LLaMA-VID	$51.25_{\pm 0.09}$	$28.28_{\pm 0.04}$	$32.85_{\pm 0.03}$	$56.59_{\pm 0.04}$	$41.22_{\pm 0.02}$	$57.49_{\pm 0.03}$	$40.39_{\pm 0.04}$
490	mPLUG-Owl	$52.73_{\pm 0.13}$	$27.47_{\pm 0.17}$	$32.47_{\pm 0.19}$	$57.60_{\pm 0.23}$	$41.32_{\pm 0.04}$	$56.32_{\pm 0.26}$	$40.83_{\pm 0.07}$
496	Chat-UniVi	$53.08_{\pm 0.01}$	$28.89_{\pm 0.02}$	$33.23_{\pm 0.08}$	$57.00_{\pm 0.06}$	$42.25_{\pm 0.04}$	$57.50_{\pm 0.03}$	$42.43_{\pm 0.03}$
407	PCC score	—	0.887	0.857	0.911	0.940	0.913	0.942
497		M3	-W3	M3-	-W4	M3-	-W5	
498	Model	L1	L2	L1	L2	L1	L2	M-avg
400	Qwen-Audio	$39.46_{\pm 0.28}$	$30.65_{\pm 0.01}$	$36.64_{\pm 0.03}$	$27.33_{\pm 0.01}$	$35.89_{\pm 0.08}$	$29.66_{\pm 0.01}$	31.84
499	Otter	$51.03_{\pm 0.04}$	$37.12_{\pm 0.00}$	$47.54_{\pm 0.00}$	$34.77_{\pm 0.00}$	$50.51_{\pm 0.03}$	$35.54_{\pm 0.00}$	39.41
500	Video-LLaMA	$47.50_{\pm 0.20}$	$36.50_{\pm 0.25}$	$52.97_{\pm 0.09}$	$35.78_{\pm 0.14}$	$46.39_{\pm 0.12}$	$34.77_{\pm 0.23}$	40.31
	VideoChat	$46.78_{\pm 0.11}$	$34.37_{\pm 0.03}$	$49.53_{\pm 0.15}$	$32.82_{\pm 0.01}$	$45.93_{\pm 0.18}$	$32.85_{\pm 0.04}$	37.58
501	SECap	$50.77_{\pm 0.03}$	$40.49_{\pm 0.03}$	$50.43_{\pm 0.03}$	$38.21_{\pm 0.03}$	$49.97_{\pm 0.03}$	$40.25_{\pm 0.03}$	42.43
502	Video-LLaVA	$52.29_{\pm 0.05}$	$40.58_{\pm 0.15}$	$52.45_{\pm 0.06}$	$39.91_{\pm 0.13}$	$52.97_{\pm 0.10}$	$39.69_{\pm 0.10}$	43.27
502	SALMONN	$56.25_{\pm 0.01}$	$43.01_{\pm 0.02}$	$50.53_{\pm 0.09}$	$38.54_{\pm 0.03}$	$53.65_{\pm 0.04}$	$42.09_{\pm 0.02}$	43.53
503	VideoChat2	$52.14_{\pm 0.23}$	$40.57_{\pm 0.14}$	$50.63_{\pm 0.19}$	$39.64_{\pm 0.18}$	$51.37_{\pm 0.14}$	$39.89_{\pm 0.15}$	42.65
504	Video-ChatGPT	$55.50_{\pm 0.13}$	$44.15_{\pm 0.18}$	$55.24_{\pm 0.02}$	$42.42_{\pm 0.05}$	$52.93_{\pm 0.05}$	$41.54_{\pm 0.14}$	46.02
304	LLaMA-VID	$55.12_{\pm 0.05}$	$44.06_{\pm 0.01}$	$56.62_{\pm 0.15}$	$42.42_{\pm 0.03}$	$53.03_{\pm 0.08}$	$41.65_{\pm 0.04}$	45.81
505	mPLUG-Owl	$55.67_{\pm 0.19}$	$43.71_{\pm 0.13}$	$55.06_{\pm 0.17}$	$40.67_{\pm 0.19}$	$54.44_{\pm 0.13}$	$42.00_{\pm 0.18}$	45.63
= 0.0	Chat-UniVi	$56.80_{\pm 0.01}$	$45.66_{\pm 0.05}$	$55.86_{\pm 0.07}$	$41.97_{\pm 0.09}$	$55.81_{\pm 0.02}$	$43.61_{\pm 0.05}$	46.75
506	PCC score	0.904	0.927	0.899	0.922	0.885	0.894	0.942

7 Table 4: GPT-based vs. EW-based grouping. We calculate the PCC score to reveal their correlation.

GPT-based vs. EW-based Grouping. This paper proposes two grouping strategies: GPT-based grouping and EW-based grouping. In this section, we explore the relationship between them and analyze whether the EW-based method can replace the GPT-based method. We calculate the scores for different types of EW-based grouping strategies, with the average score denoted as M-avg. Table 4 reports F_s , as this metric is used for the final ranking, and we also compute the PCC scores between different metrics. We observe that the PCC score between M-avg and GPT is relatively high, indicating that EW-based metrics can serve as alternatives to GPT-based metrics.

515 516 517

507

486

6 LIMITATIONS

518 Firstly, the main contribution of this paper is the definition of a new task and the conduct of foun-519 dational research. In the future, we plan to design more effective frameworks to solve OV-MER. 520 Specifically, we plan to incorporate more emotion-related instruction datasets to finetune MLLMs, 521 thereby enhancing their emotion recognition ability. Meanwhile, as mentioned in our paper, how to 522 integrate subtitle information and fuse multimodal inputs also plays a crucial role in the final perfor-523 mance. We will also consider these aspects in the framework design. Secondly, we have evaluated 524 some representative MLLMs, but not all models are covered. In the future, we will expand the scope 525 of evaluation to cover more emerging MLLMs to enrich the benchmark. Thirdly, this paper does 526 not involve cultural differences. Specifically, our original data is in Chinese, and the annotators we hired are also native Chinese speakers. In the future, we will also try to extend our method to other 527 cultures and further analyze cultural differences. 528

529 530

531

7 CONCLUSION

This paper extends traditional MER to OV-MER, allowing for the prediction of arbitrary numbers
and types of emotions. To facilitate further research, we construct an initial dataset, define evaluation
metrics, and propose solutions. We observe that current MLLMs struggle to achieve satisfactory
results, as this task requires consideration of multimodal clues and subtle temporal changes, placing
higher demands on MLLMs. Additionally, EW-based metrics can replace GPT-based metrics, thus
reducing evaluation costs while ensuring reproducibility. This paper advances current research from
basic to nuanced emotion recognition, which is crucial for building emotion AI.

540 **ETHICS STATEMENT**

541

542 The raw data of the OV-MERD dataset comes from the MER2023 dataset, from which we evenly 543 select some samples with further annotation. Therefore, we do not collect new data; we just re-544 annotate existing data. This annotation process has received consent from the dataset owners and 545 has passed our internal review. During the annotation process, we generously pay each annotator 546 approximately ¥3,000 (around \$280), which is considered high. After proofreading our annotation results, we find that the annotations focus on the multimodal clues present in the videos, without 547 any discriminatory annotations. Additionally, we restrict the use of the OV-MERD dataset to non-548 commercial purposes under the CC BY-NC 4.0 license. This license clearly outlines the correct and 549 responsible use of our dataset. Details of our license are provided in the supplementary materials. 550 In summary, this paper does not involve any ethical issues. 551

552 553

554 555

556

557 558

560

561

576

REPRODUCIBILITY STATEMENT

This paper provides the source code and intermediate results in the supplementary materials. Since we build a dataset for OV-MER, we also offer a complete description of the data construction process in Section 2 and Appendix. In summary, we have made every effort to ensure its reproducibility.

559 REFERENCES

- Felix Burkhardt, Astrid Paeschke, Miriam Rolfes, Walter F Sendlmeier, Benjamin Weiss, et al. A database of german emotional speech. In *Interspeech*, volume 5, pp. 1517–1520, 2005. 562
- 563 Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jean-564 nette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic 565 motion capture database. Language Resources and Evaluation, 42:335–359, 2008.
- 566 Carlos Busso, Srinivas Parthasarathy, Alec Burmania, Mohammed AbdelWahab, Najmeh Sadoughi, 567 and Emily Mower Provost. Msp-improv: An acted corpus of dyadic interactions to study emotion 568 perception. IEEE Transactions on Affective Computing, 8(1):67-80, 2016. 569
- Houwei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini 570 Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. IEEE Transactions on 571 Affective Computing, 5(4):377–390, 2014. 572
- 573 Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and 574 Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale 575 audio-language models. arXiv preprint arXiv:2311.07919, 2023.
- Giovanni Costantini, Iacopo Iaderola, Andrea Paoloni, Massimiliano Todisco, et al. Emovo corpus: 577 an italian emotional speech database. In Proceedings of the ninth international conference on 578 language resources and evaluation (LREC'14), pp. 3501–3504. European Language Resources 579 Association (ELRA), 2014. 580
- 581 Charles Darwin. The Expression of Emotions in Man and Animals. Penguin Classics, 1872.
- 582 Abhinav Dhall, OV Ramana Murthy, Roland Goecke, Jyoti Joshi, and Tom Gedeon. Video and 583 image based emotion recognition challenges in the wild: Emotiw 2015. In Proceedings of the 584 2015 ACM on international conference on multimodal interaction, pp. 423-426, 2015. 585
- Abhinav Dhall, Roland Goecke, Shreya Ghosh, Jyoti Joshi, Jesse Hoey, and Tom Gedeon. From 586 individual to group-level emotion recognition: Emotiw 5.0. In Proceedings of the 19th ACM 587 *international conference on multimodal interaction*, pp. 524–528, 2017. 588
- 589 Mathilde M Duville, Luz M Alonso-Valerdi, and David I Ibarra-Zarate. The mexican emotional 590 speech database (mesd): elaboration and assessment based on machine learning. In 2021 43rd An-591 nual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC), 592 pp. 1644-1647. IEEE, 2021.

Paul Ekman. An argument for basic emotions. Cognition & emotion, 6(3-4):169-200, 1992.

607

621

631

- ⁵⁹⁴ C Fabian Benitez-Quiroz, Ramprakash Srinivasan, and Aleix M Martinez. Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5562–5570, 2016.
- Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation
 with image-level labels. In *European Conference on Computer Vision*, pp. 540–557. Springer, 2022.
- Ian J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner,
 Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. Challenges in representation
 learning: A report on three machine learning contests. In *Neural information processing: 20th international conference, ICONIP 2013, daegu, korea, november 3-7, 2013. Proceedings, Part III* 20, pp. 117–124. Springer, 2013.
- Michael Grimm, Kristian Kroschel, and Shrikanth Narayanan. The vera am mittag german audio visual emotional speech database. In *IEEE International Conference on Multimedia and Expo*,
 pp. 865–868. IEEE, 2008.
- Jiaming Han, Kaixiong Gong, Yiyuan Zhang, Jiaqi Wang, Kaipeng Zhang, Dahua Lin, Yu Qiao, Peng Gao, and Xiangyu Yue. Onellm: One framework to align all modalities with language. *arXiv preprint arXiv:2312.03700*, 2023.
- Philip Jackson and SJUoSG Haq. Surrey audio-visual expressed emotion (savee) database. University of Surrey: Guildford, UK, 2014.
- Jesin James, Li Tian, and Catherine Watson. An open source emotional speech corpus for human robot interaction applications. *Interspeech 2018*, 2018.
- ⁶²⁰ William James. What is emotion? *Mind*, 9(34):188—-205, 1884.
- Kingxun 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*, pp. 2881–2889, 2020.
- Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual representation empowers large language models with image and video understanding. *arXiv* preprint arXiv:2311.08046, 2023.
- 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.
- Jean Kossaifi, Robert Walecki, Yannis Panagakis, Jie Shen, Maximilian Schmitt, Fabien Ringeval, Jing Han, Vedhas Pandit, Antoine Toisoul, Björn Schuller, et al. Sewa db: A rich database for audio-visual emotion and sentiment research in the wild. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 43(3):1022–1040, 2019.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A
 multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023a.
- Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven
 semantic segmentation. In *International Conference on Learning Representations*, 2021.
- KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023b.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,
 Ping Luo, Limin Wang, and Yu Qiao. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

660

661 662

675

683

688

689

690

- Shan Li, Weihong Deng, and JunPing Du. Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2852–2861, 2017.
- Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language
 models. *arXiv preprint arXiv:2311.17043*, 2023c.
- ⁶⁵⁴
 ⁶⁵⁵
 ⁶⁵⁶
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁴
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁴
 ⁶⁵⁴
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵¹
 ⁶⁵²
 ⁶⁵³
 ⁶⁵⁵
 ⁶⁵⁵
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵¹
 ⁶⁵²
 ⁶⁵³
 ⁶⁵⁴
 ⁶⁵⁵
 ⁶⁵⁵
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵¹
 ⁶⁵²
 ⁶⁵³
 ⁶⁵⁵
 ⁶⁵⁵
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁶
 ⁶⁵⁷
 ⁶⁵⁷
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁸
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 ⁶⁵⁹
 <
 - Zheng Lian, Haiyang Sun, Licai Sun, Zhuofan Wen, Siyuan Zhang, Shun Chen, Hao Gu, Jinming Zhao, Ziyang Ma, Xie Chen, et al. Mer 2024: Semi-supervised learning, noise robustness, and open-vocabulary multimodal emotion recognition. *arXiv preprint arXiv:2404.17113*, 2024a.
- Zheng Lian, Licai Sun, Yong Ren, Hao Gu, Haiyang Sun, Lan Chen, Bin Liu, and Jianhua Tao.
 Merbench: A unified evaluation benchmark for multimodal emotion recognition. *arXiv preprint arXiv:2401.03429*, 2024b.
- Zheng Lian, Licai Sun, Haiyang Sun, Kang Chen, Zhuofan Wen, Hao Gu, Bin Liu, and Jianhua Tao.
 Gpt-4v with emotion: A zero-shot benchmark for generalized emotion recognition. *Information Fusion*, 108:102367, 2024c.
- Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pp. 740–755. Springer, 2014.
- Rui Liu, Haolin Zuo, Zheng Lian, Xiaofen Xing, Björn W Schuller, and Haizhou Li. Emotion and intent joint understanding in multimodal conversation: A benchmarking dataset. *arXiv preprint arXiv:2407.02751*, 2024.
- Yihe Liu, Ziqi Yuan, Huisheng Mao, Zhiyun Liang, Wanqiuyue Yang, Yuanzhe Qiu, Tie Cheng, Xiaoteng Li, Hua Xu, and Kai Gao. Make acoustic and visual cues matter: Ch-sims v2. 0 dataset and av-mixup consistent module. In *Proceedings of the International Conference on Multimodal Interaction*, pp. 247–258, 2022a.
- 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*, pp. 24–32, 2022b.
 - 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.
- Reza Lotfian and Carlos Busso. Building naturalistic emotionally balanced speech corpus by retriev ing emotional speech from existing podcast recordings. *IEEE Transactions on Affective Comput- ing*, 10(4):471–483, 2017.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts.
 Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, 2011.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pp. 12585–12602, 2024.

702 Olivier Martin, Irene Kotsia, Benoit Macq, and Ioannis Pitas. The enterface'05 audio-visual emotion 703 database. In Proceedings of the 22nd International Conference on Data Engineering Workshops, 704 pp. 1-8. IEEE, 2006. 705 Gary McKeown, Michel Valstar, Roddy Cowie, Maja Pantic, and Marc Schroder. The semaine 706 database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent. *IEEE Transactions on Affective Computing*, 3(1):5–17, 2011. 708 709 Marvin Minsky. Society of mind. Simon and Schuster, 1988. 710 711 Ali Mollahosseini, Behzad Hasani, and Mohammad H Mahoor. Affectnet: A database for facial expression, valence, and arousal computing in the wild. IEEE Transactions on Affective Computing, 712 10(1):18-31, 2017. 713 714 Louis-Philippe Morency, Rada Mihalcea, and Payal Doshi. Towards multimodal sentiment analy-715 sis: Harvesting opinions from the web. In Proceedings of the 13th International Conference on 716 Multimodal Interfaces, pp. 169–176, 2011. 717 Gpt-4v(ision) system card, 2023. URL https://openai.com/research/ 718 OpenAI. 719 gpt-4v-system-card. 720 Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Thumbs up? sentiment classification using 721 machine learning techniques. In Proceedings of the Conference on Empirical Methods in Natural 722 Language Processing, pp. 79-86, 2002. 723 724 Verónica Pérez-Rosas, Rada Mihalcea, and Louis-Philippe Morency. Utterance-level multimodal 725 sentiment analysis. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, pp. 973–982, 2013. 726 727 Robert Plutchik. A general psychoevolutionary theory of emotion. *Emotion: Theory, research, and* 728 experience, 1, 1980. 729 730 Robert Plutchik. The nature of emotions: Human emotions have deep evolutionary roots, a fact that 731 may explain their complexity and provide tools for clinical practice. American Scientist, 89(4): 344-350, 2001. 732 733 Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada 734 Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations. In 735 Proceedings of the 57th Conference of the Association for Computational Linguistics, pp. 527– 736 536, 2019. 737 738 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, 739 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In Proceedings of the Conference on Empirical Methods in Natural Language Process-740 *ing*, pp. 1631–1642, 2013. 741 742 Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model 743 to instruction-follow them all. In Proceedings of the 1st Workshop on Taming Large Language 744 *Models: Controllability in the era of Interactive Assistants*, pp. 11–23, 2023. 745 746 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, 747 2024. 748 749 Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, 750 and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. In 751 Proceedings of the Twelfth International Conference on Learning Representations, 2023. 752 Maxim Tkachenko, Mikhail Malyuk, Andrey Holmanyuk, and Nikolai Liubimov. 753 Label Studio: Data labeling software, 2020. URL https://github.com/heartexlabs/ 754 label-studio. Open source software available from https://github.com/heartexlabs/label-755

studio.

- Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pp. 6558– 6569, 2019.
- Gyanendra K Verma and Uma Shanker Tiwary. Affect representation and recognition in 3d continuous valence–arousal–dominance space. *Multimedia Tools and Applications*, 76:2159–2183, 2017.
- Amy Beth Warriner, Victor Kuperman, and Marc Brysbaert. Norms of valence, arousal, and dominance for 13,915 english lemmas. *Behavior research methods*, 45:1191–1207, 2013.
- Martin Wöllmer, Felix Weninger, Tobias Knaup, Björn Schuller, Congkai Sun, Kenji Sagae, and Louis-Philippe Morency. Youtube movie reviews: Sentiment analysis in an audio-visual context. *IEEE Intelligent Systems*, 28(3):46–53, 2013.
- Chung-Hsien Wu, Jen-Chun Lin, and Wen-Li Wei. Survey on audiovisual emotion recognition:
 databases, features, and data fusion strategies. *APSIPA Transactions on Signal and Information Processing*, 3, 2014.
- Jianzong Wu, Xiangtai Li, Shilin Xu, Haobo Yuan, Henghui Ding, Yibo Yang, Xia Li, Jiangning Zhang, Yunhai Tong, Xudong Jiang, et al. Towards open vocabulary learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Yaoxun Xu, Hangting Chen, Jianwei Yu, Qiaochu Huang, Zhiyong Wu, Shi-Xiong Zhang, Guangzhi Li, Yi Luo, and Rongzhi Gu. Secap: Speech emotion captioning with large language model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 19323–19331, 2024.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- Wenmeng Yu, Hua Xu, Fanyang Meng, Yilin Zhu, Yixiao Ma, Jiele Wu, Jiyun Zou, and Kaicheng Yang. Ch-sims: A chinese multimodal sentiment analysis dataset with fine-grained annotation of modality. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 3718–3727, 2020.
- Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Tensor
 fusion network for multimodal sentiment analysis. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1103–1114, 2017.
- 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 Linguis- tics (Volume 1: Long Papers)*, pp. 2236–2246, 2018.
- Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14393–14402, 2021.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language
 model for video understanding. In *Proceedings of the Conference on Empirical Methods in Nat- ural Language Processing: System Demonstrations*, pp. 543–553, 2023.
- Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. From facial expression recognition to interpersonal relation prediction. *International Journal of Computer Vision*, 126:550–569, 2018.
- 805 806
- 807
- 808
- 809

810 DETAILED MOTIVATION А

811 812

Limitations of Dimensional Models. Although dimensional models can be used to distinguish 813 different emotions, they are abstract and difficult for the general public to understand. For instance, 814 the VAD model is one of the most widely used dimensional models. Given a VAD score (V=3.3, 815 A=2.9, D=4.7), it is challenging for people to accurately interpret the corresponding emotion. Ac-816 cording to previous research (Verma & Tiwary, 2017), this VAD score represents sad. Furthermore, 817 the VAD model attempts to simplify emotions into three linear dimensions. However, emotions are 818 inherently complex and often involve combinations of multiple emotional states, which cannot be fully captured by these three dimensions. Additionally, discrete models align more closely with 819 how humans intuitively understand emotions, making it easier to achieve consensus across different 820 annotators. Therefore, this paper uses discrete models to describe emotions.

821 822

823

830

824 OV-MER vs. Existing Emotion Models. OV-MER combines the advantages of both discrete and 825 dimensional models. Compared to one-hot discrete models, OV-MER uses denser discrete labels to 826 capture more subtle emotions, mimicking the continuous attributes of dimensional models. Com-827 pared to dimensional models, OV-MER uses discrete labels, making it closer to the way humans understand emotions. Therefore, our OV-MER provides a bridge connecting discrete and dimen-828 sional models. 829

Video Emotion vs. Facial Emotion. Video emotion is more complex than facial emotion. This 831 is because, in videos, we need to capture subtle changes in the temporal dimension and integrate 832 multimodal clues. Take Figure 1(b) as an example. In the temporal dimension, we need to infer 833 a person's nervousness based on his stuttering; in the multimodal dimension, we need to combine 834 information from different modalities to gain a more comprehensive understanding of emotion. Due 835 to the complexity of video emotion, using a single label is limiting, and more discrete labels are 836 required to better describe video emotion. This is also the motivation behind our OV-MER task. 837

838 Label Importance. In OV-MER, we do not assign different levels of importance to each label. Ev-839 ery emotion holds equal significance, and neglecting anyone can impact the performance of down-840 stream tasks. For example, if a human-computer interaction system only captures basic emotions 841 while overlooking nuanced ones, it may fail to generate appropriate responses. 842

843 844

845

В **RELATED WORK**

Multimodal Emotion Recognition. MER has rapidly developed in recent years (Wu et al., 2014). 846 Current research mainly focuses on building more efficient architectures to achieve higher accuracy 847 on benchmark datasets (Sun et al., 2024). For example, Zadeh et al. (2017) proposed a tensor fu-848 sion network that addressed the MER task by leveraging interactions among unimodal, bimodal, 849 and trimodal inputs. Tsai et al. (2019) introduced a Transformer-based model that learned implicit 850 alignment between different modalities and achieved promising results. Lian et al. (2024b) fur-851 ther established MERBench, involving various features, fusion strategies, and datasets. In emotion 852 recognition, benchmark datasets usually limit the label space to basic emotions and use majority 853 voting to determine the most likely one or more labels (Lian et al., 2023; Li et al., 2017). However, 854 emotional categories extend far beyond basic emotions. Restricting the label space will inevitably 855 overlook some nuanced emotions. To address this issue, we extend traditional MER to OV-MER, 856 which allows for the prediction of any number and categories of emotions.

857

858 Open Vocabulary Learning. Its main goal is to identify categories beyond the annotated label 859 space (Wu et al., 2024), which has been applied in various fields, such as object detection (Zareian 860 et al., 2021), segmentation (Ghiasi et al., 2022), and scene understanding (Li et al., 2021). For exam-861 ple, the object detection dataset COCO (Lin et al., 2014) contains 80 categories, while objects in the real world are nearly infinite, highlighting the importance of open vocabulary learning. This paper 862 makes the first attempt to address MER in an open-vocabulary manner. Compared to other tasks, MER is more difficult as it requires considering multimodal clues and subtle temporal variations.

C MORE EXAMPLES

In Figures 9~11, we provide more examples to visualize the difference between one-hot and OV labels. This paper uses emotion-related descriptions as a bridge for OV label extraction. Since the original video contains real people, we use DemoAI to remove personal information to address copyright concerns. Our OV-MERD dataset is derived from the MER2023 dataset (Lian et al., 2023) with further annotations. Therefore, for the original data, please download the MER2023 dataset.



Subtitle: Goodness, the relationships on our 22nd floor are getting more and more complicated. It seems like I'll have to speak less in the future.

One-hot Label: surprise

Description: In the video, there are two ladies in the frame. The lady on the left has her back to the camera, while the lady on the right is facing the camera. We mainly analyze the emotional state of the lady on the right. At the beginning of the video, her facial expression is not positive, and she seems a bit nervous and uneasy. In the following scenes, her eyebrows are raised, indicating that she is emphasizing something. At the end of the video, she lowers her head and looks down, not looking at the other person while speaking. This behavior may stem from inner pressure, anxiety, uneasiness, or frustration, causing her to tend to hide her emotions or avoid communication with others. In the audio, the volume is low, the tone is low, and there is a sense of pressure in the voice. In the text, the subtitle says, "Goodness, the relationships on our 22nd floor are getting more and more complicated. It seems like I'll have to speak less in the future." Based on the tense and uneasy facial expression, the emphasized movement of raising eyebrows, and the behavior of lowering her head and looking down, it can be inferred that she may be feeling pressure, anxiety, uneasiness, or frustration. Additionally, based on the audio clue describing low volume, low tone, and a sense of pressure in the voice, it can be further confirmed that her emotional state is likely negative. Therefore, this sentence expresses her exclamation about the increasing complexity of the relationship on the 22nd floor and her decision to speak less in the future, implying her uneasiness and pressure. Description \rightarrow OV labels: pressure, negative, uneasy, surprised

Figure 9: Example1.

918	
919	9 R. 4 S H
920	
921	
922	
923	Constitut
924	Subtitle: Reading requires literacy! In movies, even the illiterate can enjoy themselves, but what
925	really competes with you is cross-talk.
926	One-hot Label: <i>happy</i>
927	Description: In the video, at the beginning, we see a woman sitting indoors with bright
928	windows and green plants in the background. Her facial expression is calm, with her head tilted
929	to the right and her mouth slightly open, as if she is engaged in a lighthearted conversation. As
930	time goes by, her smile becomes more pronounced, and she slightly shakes her head to the left,
931	glancing to the left as if referring to someone or something out of view. Towards the end of the
932	something to the other person. In the audio, the character's tone is relatively steady, laced with
933	humor. However, in the first sentence, there are noticeable emphases: "reading" has an upward
934	intonation, while "literacy" has a downward intonation, creating a contrast. The subtitles read:
935	"Reading requires literacy! In movies, even the illiterate can enjoy themselves, but what really
936	competes with you is cross-talk." This sentence might be the woman's comment or reaction to
937	someone hearby. Given the autio cues of a steady tone with humor and the woman's sinne and glance to the left, we can infer that her words carry a mocking or sarcastic tone. Thus, this
938	sentence likely expresses the woman's mockery of the other person's lack of understanding or
939	recognition of something, conveyed in a humorous manner.
940	Description \rightarrow OV labels: <i>humor, mockery, irony</i>
041	
040	Figure 10: Example2.
942	
943	
944	
945	
946	
947	
948	
949	Constant River River Constant Con
950	Subtitle: Luyuan will never fall into your hands, because he has found a high-paying job.
951	One-hot Label: angry
952	Description: In the video, the frame includes two women. The woman on the left has her back
953	to the camera, while the woman on the right is facing the camera. We mainly analyze the
954	emotional state of the woman on the right. At the beginning of the video, she slightly furrows
955	her brows, opens her mouth, and her facial expression appears angry, as if she is engaged in a
956	heated argument with the other person. At the same time, she points her finger at the other person and makes a motion towards them, which may indicate that she is accusing the other
957	person or emphasizing her own viewpoint. Overall, she may be going through a debate or
958	intense conversation, and her emotional state may be one of excitement and anger. She seems to
959	be accusing and expressing her dissatisfaction. In the audio, the tone is aggressive and the
960	character's emotions are more excited. Combined with the text content, the tone seems to carry a
961	sense of threat. In the text, the subtile says, "Luyuan will never fail into your hands, because he
962	the right to the woman on the left. Based on the angry and angry emotions displayed by the
963	woman on the right in the video clues, as well as her pointing finger and motion towards the
964	other person, it can be inferred that she is accusing the other person or emphasizing her own
965	viewpoint. At the same time, based on the aggressive tone and excited emotions described in the
966	can be inferred that this sentence may carry a sense of threat and the woman on the right may be
967	threatening the woman on the left not to interfere with or harm Luyuan. Therefore, this sentence
968	expresses the woman on the right's anger and threatening emotions.
969	Description \rightarrow OV labels: <i>warning, angry, threat</i>
070	
970	Figure 11: Example3.
311	

972 D SUMMARY OF ABBREVIATIONS

In Table 5, we summarize the main abbreviations and their meanings.

Category	Abbreviation	Explanation
T 1 1	OH label	One-hot label, the most likely label in a limited set of basic en
Label	OV labels	Open-vocabulary labels, a set of labels in an unlimited label s
	MER	Multimodal Emotion Recognition, which aims to recognize th emotion label.
Task	OV-MER	Open-vocabulary MER, which aims to identify the OV emotion
Dataset	OV-MERD	This is a dataset we built for the OV-MER task.
	EW	Emotion Wheel.
Metric	M1, M2, M3	Different grouping strategies based on the emotion wheel.
	Precision _S , Recall _S , F _S	Metrics defined for OV-MER.
	LLM	Large Language Model. Large-scale models and only process
Madal	ALLM	Audio LLM. Different from LLM, it can also process audio ir
Woder	VLLM	Video LLM. Different from LLM, it can also process video in
	MLLM	Multimodal LLM. Unlike LLM, it can process at least one monity (e.g., audio or video). Thus, MLLM includes ALLM and V
	CLUE-Multi	It uses the checked acoustic and visual clues to generate descr
	CLUE-Audio	Different from CLUE-Multi, it only uses checked acoustic clu
	CLUE-Video	Different from CLUE-Multi, it only uses checked visual clues
Description	CLUE-Text	It only relies on text to generate descriptions.
Description	CLUE-A/T/V	Any of CLUE-Audio, CLUE-Text, and CLUE-Video.
	CLUE-M/A/T/V	Any of CLUE-Multi, CLUE-Audio, CLUE-Text, and CLUE-
	CLUE-MLLM	It uses the output from MLLM without any manual checking
	S0, S1, S2	Different CLUE-MLLM generation strategies.

Е DETAILS IN DATASET CONSTRUCTION

Prompts. Figure 2 presents our dataset construction process. In Table 6, we provide prompts and corresponding models used in this process.

Table 6: Prompts and corresponding models used in the dataset construction process.

1033	Function (Model)	Prompt
1024		As an expert in the field of emotions, please focus on facial expressions, body language,
1034		environmental cues, and events in the video and predict the emotional state of the character.
1035	#1 Pre-label visual clue	Please ignore the character's identity. We uniformly sample 3 frames from this video. Please
1036	(VLLM)	consider the temporal relationship between these frames and provide a complete description
1007		of this video. Avoid using descriptions like "the first image" and "the second image", and
1037		instead use terms like "beginning", "middle", and "end" to denote the progression of time.
1038	#2 Pre-label acoustic clue	As an expert in the field of emotions, please focus on the acoustic information in the audio to
1039	(ALLM)	discern clues related to the emotions of the individual. Please provide a detailed description
1000	(TEELIN)	and ultimately predict the emotional state of the individual.
1040		Please act as an expert in the field of emotions. We provide acoustic and visual clues that
1041	#3 Merge	may be related to the character's emotional state, along with the original subtitle of the
1040	(LLM)	video. Please analyze which parts can infer the emotional state and explain the reasons.
1042		During the analysis, please integrate the textual, audio, and visual clues.
1043	#4 Translation	<i>Chinese</i> \rightarrow <i>English:</i> Please translate the following sentence from Chinese into English.
1044	(LLM)	$English \rightarrow Chinese$: Please translate the following sentence from English into Chinese.
10/5		Please assume the role of an expert in the field of emotions. We provide clues that may
1045	#5 OV label extraction (LLM)	be related to the emotions of the characters. Based on the provided clues, please identify
1046		the emotional states of the main characters. Please separate different emotional categories
1047		with commas and output only the clearly identifiable emotional categories in a list format.
1048		It none are identified, please output an empty list.

Merging Process. In this paper, we rely on the powerful reasoning capabilities of LLM for mul-timodal fusion. Specifically, as shown in Table 6, we ask LLM to integrate textual, acoustic, and visual clues to infer the emotional state. From the output, we observe that LLM can produce reason-able analytical results. However, ambiguities and contradictions in multimodal fusion are inevitable, which is a challenging and open problem. In this paper, we simply use LLM to address this problem, providing a practical solution for the OV-MER task. However, more effective strategies may exist, such as using more complex prompts or incorporating modality control measures. Therefore, we plan to explore this area in the future.

1080 **Visualization.** In Figure 12, we visualize the output of the main steps. 1081



Subtitle: Why are you all looking at me like that? So, as long as it's a woman, does she have to have a relationship with me?

Visual Clue (Prelabeled): In the opening scene, we see a man wearing a white lab coat. He looks confident, with a slight smile on his lips and eyes directly facing the camera, as if he is explaining or narrating something. The background features computers and office equipment, suggesting that this might be an office or laboratory setting. As the scene progresses to the middle part, the man's head turns slightly to one side, and his gaze shifts sideways, seemingly interacting with someone next to him. His smile broadens, and his eyebrows lift slightly, indicating that he might be having a friendly conversation or sharing a light-hearted moment. In the next scene, the man's expression becomes more exaggerated; his mouth opens wide as if expressing surprise or emphasizing a point. This kind of facial expression and body language is typically associated with sharing surprising news or reacting strongly to something. In the final scene, he returns to facing the camera directly, his expression calm once again with a slight smile, giving the impression that he might have just finished explaining or narrating something and feels satisfied or confident. Overall, based on these scenes, we can infer that the man in the video is likely conducting some sort of demonstration or explanation. He interacts with the audience, shares information, and possibly expresses surprise at some point. His emotional state throughout the video appears to be positive, exhibiting confidence and a friendly demeanor.

1104 Visual Clue (Checked): In the video, the screen shows a male character in an indoor setting. At the 1105 beginning of the video, his eyes are wide open and his mouth is also open, indicating a surprised facial 1106 expression. In the following scenes, he looks around, seemingly explaining or narrating something to 1107 the people around him. Overall, his emotions are not positive or optimistic.

1108 Acoustic Clue (Prelabeled): The character in the voice is expressing frustration and annoyance. The 1109 tone is sarcastic and mocking, and the intonation is rising at the end of each sentence. There are pauses 1110 between each sentence, and the stress is on certain words to emphasize the frustration and annoyance.

1111 Acoustic Clue (Checked): In the audio, the character speaks with a stutter, which usually expresses 1112 feelings of nervousness, anxiety, or unease. Combined with the text content, the character seems to be unhappy and angry due to the prejudice of the people around him. 1113

1114 CLUE-Multi: In the video, the screen shows a male character in an indoor setting. At the beginning of the video, his eyes are wide open and his mouth is also open, indicating a surprised facial expression. 1115 In the following scenes, he looks around, seemingly explaining or narrating something to the people 1116 around him. Overall, his emotions are not positive or optimistic. In the audio, the character speaks with 1117 a stutter, which usually expresses feelings of nervousness, anxiety, or unease. Combined with the text 1118 content, the character seems to be unhappy and angry due to the prejudice of the people around him. 1119 The subtitle in the text says, "Why are you all looking at me like that? So, as long as it's a woman, does 1120 she have to have a relationship with me?" This sentence expresses the male character's dissatisfaction and anger towards the people around him. Based on the surprised and negative facial expression of the 1122 male character in the video clues, as well as the stuttering speech in the audio clues, we can infer that 1123 the male character is expressing a feeling of dissatisfaction and anger in this sentence. He may feel 1124 troubled by the prejudice of the people around him and is unhappy with this unfair treatment.

Figure 12: An example to visualize the output of the main steps.

1126 1127 1128

1125

1121

1082 1083 1084

1087 1088

1089

1090

1093

1094

1095

1099

1100

1101

1102

1103

- 1129
- 1130
- 1131

1132

¹¹³⁴ F ONE-HOT VS. OV LABELS

1155

1156

This section provides a deeper comparison between the one-hot labels in the MER2023 dataset and the OV labels in the OV-MERD dataset. Figure 13 shows the word cloud and label number distribution of OV labels. In Figure 13(a), we observe that OV labels cover a wider variety of emotions, some of which (such as *shy*, *nervous*, and *grateful*) are rarely discussed in previous datasets. In Figure 13(b), we notice that most samples have about 2 to 4 labels, much more than the traditional task where each sample is assigned only one emotion. Therefore, OV-MER provides richer labels.





1157 In Table 7, we report the performance of one-hot labels in OV-MER. We observe that one-hot labels 1158 have high *precisions* but low *recalls*, indicating that one-hot labels are correct but not comprehen-1159 sive. Due to the limited label space and the constrained number of labels, one-hot labels cannot 1160 cover all emotions, highlighting the limitations of traditional MER and the importance of OV-MER. 1161 Additionally, these results reflect the necessity to use F_s for the final ranking, which can balance 1162 accuracy and completeness.

Table 7: Performance of one-hot labels in OV-MER.						
Language	$F_{S}\uparrow$	Precision _s \uparrow	$\text{Recall}_{s} \uparrow$			
English	$65.71_{\pm 0.06}$	$92.17_{\pm 0.00}$	$51.05_{\pm 0.08}$			
Chinese	$66.16_{\pm 0.02}$	$93.07_{\pm 0.00}$	$51.32_{\pm 0.03}$			

1169Figure 14 shows the emotion distribution of OV labels. We observe that the number of samples for1170different emotions follows a long-tail distribution. These results indicate that OV labels not only1171cover some major labels but also capture subtle emotions that occur infrequently.



1188 DATASET COMPARISON G

1189

1190 This paper introduces a new task, OV-MER, and constructs a dataset for this task called OV-MERD. 1191 Table 8 compares OV-MERD with existing datasets. The annotation types of these datasets can 1192 be broadly categorized into two types: dimensional emotions and discrete emotions. We classify 1193 sentiment analysis datasets (e.g., CMU-MOSI) as dimensional datasets because the definition of sentiment intensity overlaps with the valence in dimensional emotions. We observe that OV-MERD 1194 contains 248 emotion categories, with most samples having 2 to 4 labels, significantly exceeding the 1195 number of labels in existing datasets. In the future, as the scale of the dataset increases, the number 1196 of candidate labels can be further expanded. Meanwhile, we would like to emphasize that our 1197 OV-MERD is the first dataset that uses the human-LLM collaborative annotation strategy, aiming 1198 to provide richer labels to capture more nuanced emotions. We believe this work is an important 1199 extension of traditional MER and will contribute to the development of the field. 1200

1201

1202 Table 8: Dataset comparison. In this table, "I", "A", "V", and "T" are abbreviations for image, audio, 1203 video, and text, respectively. Some datasets (such as CMU-MOSEI and MSP-Podcast) contain both discrete and dimensional amotions 1004

1204	Dataset	Modality	Annotation Type	# Categories	# Labels per Sample
1205	MSP-Podcast (Lotfian & Busso, 2017)	A	Dimensional	3	1
1206	SST (Socher et al., 2013)	Т	Dimensional	1	1
1207	Cornell (Pang et al., 2002)	Т	Dimensional	1	1
1000	Large Movie (Maas et al., 2011)	Т	Dimensional	1	1
1200	ICT-MMMO (Wöllmer et al., 2013)	A,V,T	Dimensional	1	1
1209	YouTube (Morency et al., 2011)	A,V,T	Dimensional	1	1
1210	MOUD (Pérez-Rosas et al., 2013)	A,V,T	Dimensional	1	1
1011	CMU-MOSI (Zadeh et al., 2017)	A,V,T	Dimensional	1	1
1211	CMU-MOSEI (Zadeh et al., 2018)	A,V,T	Dimensional	1	1
1212	CH-SIMS (Yu et al., 2020)	A,V,T	Dimensional	1	1
1213	CH-SIMS v2 (Liu et al., 2022a)	A,V,T	Dimensional	1	1
1014	VAM (Grimm et al., 2008)	A,V,T	Dimensional	3	1
1214	SEMAINE (McKeown et al., 2011)	A,V,T	Dimensional	5	1
1215	AFEW-VA (Kossaifi et al., 2017)	A,V,T	Dimensional	2	1
1216	SEWA(Kossaifi et al., 2019)	A,V,T	Dimensional	3	1
4047	MSP-Podcast (Lotfian & Busso, 2017)	A	Discrete	8	1
1217	JL-Corpus (James et al., 2018)	A	Discrete	10	1
1218	EmoDB (Burkhardt et al., 2005)	A	Discrete	7	1
1219	EMOVO (Costantini et al., 2014)	A	Discrete	7	1
1210	MESD (Duville et al., 2021)	A	Discrete	6	l
1220	SFEW 2.0 (Dhall et al., 2015)	1	Discrete	7	
1221	FER-2013 (Goodfellow et al., 2013)	l	Discrete	/	1
1222	EmotioNet (Fabian Benitez-Quiroz et al., 2016)		Discrete	23	
4000	AffectNet (Mollahosseini et al., 2017)		Discrete	7	
1223	Expw (Zhang et al., 2018)		Discrete	/	
1224	RAF-DB(Li et al., 2017)		Discrete	19	1~2
1225	CMU-MOSEI (Zaden et al., 2018)	A,V,I	Discrete	6	1
1000	en lekface (martin et al., 2000)	A,V,I	Discrete	0	1
1226	SAVEE (Jackson & Haq, 2014) AFEW 7.0 (Dhall at al. 2017)	A,V,I	Discrete	7	1
1227	AFEW (1.0) (Ditallet al., 2017) MAEW (1.0) at al. 2022 b)	A,V,I	Discrete	/ 11	1
1228	DEFW (Ling et al. 2020)	A,V,I	Discrete	11	1
1000	CPEMA D (Coo et al. 2014)	A, V, I	Discrete	6	1
1229	MSP-IMPROV (Busso et al. 2016)	ΔVT	Discrete	0	1
1230	RAVDESS Livingstone & Russo (2018)	A VT	Discrete	4	1
1231	IEMOCAP (Busso et al., 2008)	A.V.T	Discrete	10	1
1020	MELD (Poria et al., 2019)	A.V.T	Discrete	7	1
1232	MC-EIU (Liu et al., 2024)	A.V.T	Discrete	7	1
1233	MER2023 (Lian et al., 2023)	A.V.T	Discrete	6	1
1234	MER2024 (Lian et al., 2024a)	A,V,T	Discrete	6	1
1235	OV-MERD (Ours)	A,V,T	Discrete	248	$1 \sim 9$, most $2 \sim 4$
1236				(arbitrary label)	(arottrary number)

- 1237
- 1238
- 1239

1240

¹²⁴² H DURATION DISTRIBUTION OF OV-MERD

1243 1244

1248 1249 1250

1251 1252

1253

1255 1256 1257

1259

1261 1262

1263 1264 1265

1266

In Figure 15, we analyze the duration distribution of the OV-MERD dataset. We observe that the majority of the samples have durations ranging from 1 to 4 seconds. This distribution is consistent with that of the MER2023 dataset, which was used as the original dataset for constructing OV-MERD (see Section 2.3 for details).



Figure 15: Duration distribution of the OV-MERD dataset.

I NUMBER OF SAMPLED FRAMES IN PRE-ANNOTATION

To generate pre-annotated visual clues, we sample three frames from each video and input them into GPT-4V. In this section, we discuss the rationale behind the choice of the number of sampled frames. Specifically, we categorize visual clues into two types: (1) visual clues with relatively long durations; and (2) visual clues with fast movements, such as eye movements, head movements, and micro-expressions.

1272 For the first type, since the duration of most videos is between 1 and 4 seconds (see Appendix H) 1273 and the video content is usually continuous with only minor differences between adjacent frames 1274 (see Appendix C), uniformly sampling three frames are sufficient to capture this information; For 1275 the second type, we observe that current MLLMs (including the GPT-4V used in this paper) struggle 1276 to capture these fast movements. Increasing the number of sampled frames does not address this 1277 issue. Previous research has also shown that GPT-4V cannot recognize micro-expressions (Lian et al., 2024c). To capture these fast movements, we employ multiple professional annotators to 1278 manually add this information. 1279

1280

¹²⁸¹ J ANNOTATION DETAILS

1282

1283 This section presents our annotation guidelines and the layout of the annotation platform. Our 1284 annotation process relies on the Label Studio (Tkachenko et al., 2020) toolkit. As shown in Figure 1285 2, there are two parts that require manual checking: 1) the pre-annotated acoustic and visual clues; 1286 2) the merged open-vocabulary labels. To reduce subjective bias, we hire eight annotators who 1287 are experts in affective computing and familiar with the definitions of emotions. Additionally, we conduct two rounds of checks with no overlap among annotators in each round. Specifically, in the 1288 first round, we randomly select four annotators to check the clues and labels; in the second round, 1289 we merge the clues and labels reviewed by the first four annotators and ask another four annotators 1290 to perform a second round of checks. Ultimately, we find that these checked clues and labels are 1291 well-aligned with the video content. 1292

Figure 16 shows the layout of the annotation platform used for manually checking acoustic and visual clues. During the annotation process, we use the following instructions: We provide pre-labeled acoustic and visual clues. Please manually check these clues, remove errors, and add missing information. On the annotation platform, we design an interface with a time slider, allowing annotators 1297 the manual check, helping them better annotate the details that may be missed in pre-annotation.

 1298

 1299

 1300

 1301

 1302

 1303

to start playing the video from any frame. This enables annotators to view the entire video during



Figure 16: Layout of the annotation platform used for manually checking acoustic and visual clues.

1323 Figure 17 displays the layout of the annotation platform used for manually checking emotional 1324 labels. During annotation, we use the following instructions: Please select all labels that match the character's emotional state in the "candidate emotions". If the provided candidate labels cannot 1325 perfectly describe the character's emotional state, you can also manually add new labels to the 1326 "other emotions" part. Specifically, annotators need to label two parts. First, we list all candidate 1327 labels from which annotators can choose what they believe to be the correct labels; second, when 1328 the candidate labels cannot perfectly describe the emotions, annotators can manually add additional 1329 labels in the "other emotions" part. 1330





1296

1319 1320

Figure 17: Layout of the annotation platform used for manually checking emotional labels.

To illustrate which labels are removed or kept, we provide two examples, each with labels from multiple annotators. To ensure annotation quality, we hire professional annotators who are experts in affective computing and familiar with the definition of emotions. Some of these annotators are members of our team who specialize in affective computing. In Figure 18, as the character doesn't know which doctor to see, most annotators provide labels such as confused or puzzled. Based on his tone and expression, some annotators further provide labels like anxious and serious. In Figure 19, most annotators notice his *disapprove* based on the textual content. Combining other modalities, some annotators further note his *blame* and *accuse* of what others are planning to do. From these examples, we can observe that these annotators provided relatively reliable labels. However, some annotators may focus only on the most relevant labels and overlook some details. To ensure the comprehensiveness of the annotation results, we merge the labels checked by four annotators. For example, in Figure 18, the final merged labels are troubled, focused, puzzled, anxious, worried, confused, and serious. In the next round of checks, we invite another four annotators for a second round of checks. Through this process, we can ensure that each retained label is confirmed by at least one annotator in each round, thereby ensuring the comprehensiveness and accuracy of the annotation results.



Subtitle: Huh? Which director?

A1: troubled, focused, puzzled

A2: anxious, worried, confused

A3: confused, puzzled

A4: puzzled, anxious, confused, troubled, serious

Figure 18: Example1 with labels from multiple annotators.



Subtitle: It's not right to do this, even if it's for the child's good!

A1: surprised, dislike, disapprove, unexpected, worry

A2: disapprove, serious, worry

A3: surprised, serious, amazed, shocked

A4: disapprove, serious, accuse, blame

Figure 19: Example2 with labels from multiple annotators.

1404 Κ DETAILS OF LANGUAGE IMPACT EXPERIMENTS 1405

1406 **Experimental Design.** In Figure 2, we analyze from two perspectives: 1) the impact of descriptive 1407 language (Clue-Multi), and 2) the impact of abstract language (OV labels). The Y_{EE} to Y_{EC} (or 1408 Y_{CC} to Y_{CE}) experiment aims to keep the descriptive language consistent to analyze the effect of 1409 abstract language, while the Y_{CE} to Y_{EE} (or Y_{EC} to Y_{CC}) experiment aims to keep the abstract 1410 language consistent to analyze the effect of descriptive language.

1412 Jaccard Similarity Coefficient. Figure 2 uses the Jaccard similarity coefficient to measure the similarity between two sets, which is slightly different from the evaluation metrics defined in Section 1413 3. Specifically, in Section 3, we use the following metrics: 1414

1415 1416 1417

1423

1424

1425

1432

1433 1

1411

$$\operatorname{Precision}_{S} = \frac{|\mathcal{Y} \cap \mathcal{Y}|}{|\mathcal{Y}|}, \ \operatorname{Recall}_{S} = \frac{|\mathcal{Y} \cap \mathcal{Y}|}{|\mathcal{Y}|}, \ \operatorname{F}_{S} = 2 \times \frac{\operatorname{Precision}_{S} \times \operatorname{Recall}_{S}}{\operatorname{Precision}_{S} + \operatorname{Recall}_{S}}.$$
 (5)

1418 The motivation for the above metrics is that \mathcal{Y} represents the ground truth, while $\hat{\mathcal{Y}}$ represents 1419 the prediction. However, in Figure 2, the two sets of emotions are considered equally important. 1420 Therefore, we use the Jaccard similarity coefficient to measure the similarity. This metric evaluates 1421 the similarity between two sets by comparing the size of their intersection to the size of their union: 1422

Similarity_s =
$$\frac{|\mathcal{Y} \cap \hat{\mathcal{Y}}|}{|\mathcal{Y} \cup \hat{\mathcal{Y}}|}$$
. (6)

1426 L **CLUE-MULTI ANALYSIS** 1427

1428 In this section, we further analyze the reliability and comprehensiveness of CLUE-Multi from three 1429 aspects: discrete emotion recognition, dimensional emotion recognition, and visual clue statistics. 1430 Table 9 provides prompts and models for each part of the analysis. 1431

Table 9: Prompts and corresponding models used in CLUE-Multi analysis.

1/13/	Function (Model)	Prompt
1405		Please assume the role of an expert in the emotional domain. We provide clues
1435	#1 Discrete Emotion Recognition	that may be related to the emotions of the character. Based on the provided clues,
1436	(GPT-3.5)	identify the emotional states of the main characters. We provide a set of emotional
1/137		candidates, please rank them in order of likelihood from high to low. The candidate
1437		set is {happy, angry, worried, sad, surprise, neutral}.
1438		As an expert in the emotional domain, we provide clues that may be related to the
1439		emotions of characters. Based on the provided clues, please identify the overall
1440		positive or negative emotional polarity of the main characters. The output should
1441	#2 Valence Estimation (GPT-3 5)	ative emotions, 0 indicates neutral emotions, and +5 indicates extremely neg-
1442	(011 5.5)	emotions. Larger numbers indicate more positive emotions, while smaller numbers
1443		number with two decimal places, directly outputting the numerical result without
1444		including the analysis process.
1445	#3 Visual Clue Analysis	Please assume the role of an expert in the field of emotions. We provide clues
1//6		related to the emotions of the characters in the video. Please output the facial
1440	(GPT-3.5)	movements and body gestures involved in the description, separated by commas.
1447		The output format should be in list form.

1448

1449 **Discrete Emotion Recognition.** Our dataset is based on MER2023, which provides relatively 1450 reliable one-hot labels. Therefore, we attempt to determine whether these one-hot labels can be 1451 identified from CLUE-Multi. This part of the analysis aims to verify whether CLUE-Multi can 1452 cover the traditional one-hot emotion recognition task. Experimental results indicate that the top-1 1453 and top-2 scores can reach 93.48 and 96.89, respectively. Further analysis shows that the prediction 1454 errors are primarily due to the limitations of one-hot labels. For example, in Figure 1, the character 1455 shows a compound emotional state, including surprised, nervous, and unsatisfied. However, when we rank the candidate emotions, the output is: angry, surprised, worried, neutral, sad, happy. The 1456 top-1 label is *angry*, which differs from *surprise* in MER2023, leading to a prediction error. These 1457 results reveal the limitations of traditional one-hot labels in describing emotions.

Valence Estimation. Besides discrete labels, MER2023 also provides relatively reliable valence scores. Therefore, we attempt to verify whether CLUE-Multi can used for valence estimation. Through experimental analysis, we observe that the PCC score between predictions and annotations can reach 0.88, indicating that CLUE-Multi also contains clues for dimensional emotion recognition.

Visual Clue Analysis. Following that, we attempt to analyze the diversity of visual clues in CLUEMulti. Through experimental analysis, we observe that each sample has an average of 4.95 visual
clues. Therefore, we conclude that CLUE-Multi contains a wealth of clues that can help address
discrete emotion recognition and valence estimation. Additionally, these results validate the completeness and reliability of CLUE-Multi.

M VISUALIZATION OF CLUE-M/A/T/V

¹⁴⁷¹ Figure 20 provides an example and visualizes CLUE-M/A/T/V.

a we de la facta de la seconda de	e a secondaria se a se	- CARABARR NO.

Subtitle: Why are you all looking at me like that? So, as long as it's a woman, does she have to have a relationship with me?

CLUE-Multi: In the video, the screen shows a male character in an indoor setting. At the beginning of the video, his eyes are wide open and his mouth is also open, indicating a surprised facial expression. In the following scenes, he looks around, seemingly explaining or narrating something to the people around him. Overall, his emotions are not positive or optimistic. In the audio, the character speaks with a stutter, which usually expresses feelings of nervousness, anxiety, or unease. Combined with the text content, the character seems to be unhappy and angry due to the prejudice of the people around him. The subtitle in the text says, "Why are you all looking at me like that? So, as long as it's a woman, does she have to have a relationship with me?" This sentence expresses the male character's dissatisfaction and anger towards the people around him. Based on the surprised and negative facial expression of the male character is expressing a feeling of dissatisfaction and anger in this sentence. He may feel troubled by the prejudice of the people around him and is unhappy with this unfair treatment.

CLUE-Audio: In the audio, the character speaks with a stutter, which usually expresses feelings of nervousness, anxiety, or unease. Combined with the text content, the character seems to be unhappy and angry due to the prejudice of the people around him. In the text, the subtitle content: "Why are you all looking at me like that? So, as long as it's a woman, does she have to have a relationship with me?" This sentence expresses the character's dissatisfaction and anger. The stuttering described in the audio clues, which is usually associated with emotions like nervousness, anxiety, or unease, combined with the subtitle content suggests that the character feels dissatisfied and angry. The character feels dissatisfied with the prejudices of those around them and believes that they have certain expectations or demands of them. Therefore, combining the audio clues and the subtitle content, we can infer that the character's emotional state is one of nervousness, anxiety, unease, dissatisfaction, and anger.

CLUE-Text: Subtitle: "Why are you all looking at me like that? So, as long as it's a woman, does she have to have a relationship with me?" From this piece of text, we can infer that the character's emotional state is one of dissatisfaction or anger. The reasoning is as follows: 1. "Why are you all looking at me like that?" This question expresses the character's dissatisfaction with being stared at. It implies that the character feels uncomfortable or displeased with others watching him. 2. "So, as long as it's a woman, does she have to have a relationship with me?" This statement conveys the character's frustration with others' expectations or prejudices. The character feels annoyed or angry about the assumption that there should be a special relationship between him and a woman simply because she is female. In summary, the text suggests that the character is experiencing emotions of dissatisfaction or anger.

CLUE-Video: In the video, the screen shows a male character in an indoor setting. At the beginning of the video, his eyes are wide open and his mouth is also open, indicating a surprised facial expression. In the following scenes, he looks around, seemingly explaining or narrating something to the people around him. Overall, his emotions are not positive or optimistic.

Figure 20: Visualization of CLUE-M/A/T/V.

DETAILS OF CLUE-MLLM Ν

CLUE-MLLM directly utilizes the output from MLLM without any manual checking process. Table 10 provides model cards for different MLLMs. For each MLLM, we provide two types of prompts (see Table 11): one that ignores text and another that considers text. To ensure a fair comparison, we use similar prompts for audio, video, and audio-video LLMs.

Table 10: Model cards for MLLMs.

1500	Table 10: Wodel cards for WILLWIS.				
1520	Model	Link			
1521	SECap (Xu et al., 2024)	https://github.com/thuhcsi/SECap			
1522	SALMONN (Tang et al., 2023)	https://github.com/bytedance/SALMONN			
1500	Qwen-Audio (Chu et al., 2023)	https://github.com/QwenLM/Qwen-Audio			
1523	Otter (Li et al., 2023a)	https://github.com/Luodian/Otter			
1524	OneLLM (Han et al., 2023)	https://github.com/csuhan/OneLLM			
1525	PandaGPT (Su et al., 2023)	https://github.com/yxuansu/PandaGPT			
1020	VideoChat (Li et al., 2023b)	https://github.com/OpenGVLab/Ask-Anything/tree/main/video_chat			
1526	VideoChat2 (Li et al., 2024)	https://github.com/OpenGVLab/Ask-Anything/tree/main/video_chat2			
1527	Video-LLaMA (Zhang et al., 2023)	https://github.com/DAMO-NLP-SG/Video-LLaMA			
1528	Video-LLaVA (Lin et al., 2023)	https://github.com/PKU-YuanGroup/Video-LLaVA			
1520	Video-ChatGPT (Maaz et al., 2024)	https://github.com/mbzuai-oryx/Video-ChatGPT			
1529	LLaMA-VID (Li et al., 2023c)	https://github.com/dvlab-research/LLaMA-VID			
1530	mPLUG-Owl (Ye et al., 2023)	https://github.com/X-PLUG/mPLUG-Owl			
1521	Chat-UniVi (Jin et al., 2023)	https://github.com/PKU-YuanGroup/Chat-UniVi			
1551	GPT-4V (OpenAI, 2023)	https://openai.com/			
1532					

Table 11: Prompts for extracting emotion-related descriptions using MLLMs.

1525	Table 11. I tompts for extracting emotion-related descriptions using will livis.										
1555	Model	Text	Prompt								
1536 1537 1538		w/o	As an expert in the field of emotions, please focus on the acoustic information in the audio to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately predict the emotional state of the individual.								
1539 1540 1541	Audio LLM	w/	Subtitle content of the audio: {subtitle}; As an expert in the field of emotions, please focus on the acoustic information and subtitle content in the audio to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately predict the emotional state of the individual in the audio.								
1542 1543 1544		w/o	As an expert in the field of emotions, please focus on the facial expressions, body movements, environment, etc., in the video to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately predict the emotional state of the individual in the video.								
1545 1546 1547	Video LLM	w/	Subtitle content of the video: {subtitle}; As an expert in the field of emotions, please focus on the facial expressions, body movements, environment, subtitle content, etc., in the video to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately predict the emotional state of the individual.								
1548 1549 1550		w/o	As an expert in the field of emotions, please focus on the facial expressions, body move- ments, environment, acoustic information, etc., in the video to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately pre- dict the emotional state of the individual in the video.								
1551 1552 1553 1554	Audio-Video LLM	w/	Subtitle content of the video: {subtitle}: As an expert in the field of emotions, please focus on the facial expressions, body movements, environment, acoustic information, subtitle content, etc., in the video to discern clues related to the emotions of the individual. Please provide a detailed description and ultimately predict the emotional state of the individual in the video.								

This paper tests different CLUE-MLLM generation strategies: S0, S1, and S2. Experimental results are shown in Table 12. We observe that S2 generally outperforms both S0 and S1. Therefore, we adopt S2 as the default strategy.

71	Table 12: Per	formance	compariso	n of differe	nt strategie	s for genera	ating CLUE	E-MLLM.	
	Model	Strategy	E	English	Decell	Б	Chinese	Decell	
		50	Γ _S	40.41	20.48	Γ _S	25.71	27.51	
	Ottor	SU S1	34.75 ± 0.02	40.41 ± 0.03	10.46 ± 0.01	31.06 ± 0.10	33.71 ± 0.15	$27.31_{\pm 0.07}$	
	Otter	\$2	22.34 ± 0.05	20.03 ± 0.08	19.80 ± 0.04 38.00	25.00 ± 0.04	29.14 ± 0.03	21.99 ± 0.05	
		52 50	26.00 ± 26.09	$\frac{30.71 \pm 0.10}{20.18 \pm 0.00}$	$\frac{36.09 \pm 0.09}{25.10}$	$\frac{40.22 \pm 0.01}{28.70}$	32.05 ± 0.16	$\frac{41.10 \pm 0.08}{26.76 \pm 0.08}$	
	PandaGPT	S1	20.99 ± 0.01 34.75 ± 0.01	29.10 ± 0.08	32.94 ± 0.04	34.74 ± 0.01	37.27 ± 0.00	20.70 ± 0.03	
		\$2	45.89 ± 0.21	50.03 ± 0.01	42.39 ± 0.14	$47.33_{\pm 0.04}$	57.27 ± 0.15	42.75 ± 0.18	
		52 50	34.77 ± 0.20	37.66 ± 0.01	$\frac{42.30\pm0.33}{32.30\pm0.00}$	37.62 ± 0.04	$\frac{33.01\pm0.08}{40.33\pm0.05}$	$\frac{42.75\pm0.11}{35.25\pm0.05}$	
	Video-ChatGPT Video-LLaMA VideoChat VideoChat2	S1	41.74 ± 0.04	45.59 ± 0.13	32.30 ± 0.03 38 49 + 0.03	40.81 ± 0.02	45.03 ± 0.05	37.23 ± 0.25	
	video chator i	S2	50.52 ± 0.06	54.03 ± 0.24	47.44 ± 0.23	5473 ± 0.03	61.15 ± 0.00	49.52 ± 0.05	
		50	28.17 ± 0.06	28.64 ± 0.04	27.72 ± 0.18	30.70 ± 0.11	30.09 ± 0.14	31.34 ± 0.08	
	Video-LLaMA	S1	$34.43_{\pm 0.16}$	35.82 ± 0.30	$33.15_{\pm 0.11}$	34.01 ± 0.25	$35.16_{\pm 0.22}$	32.94 ± 0.08	
	VIGEO ELaivir Y	S2	44.73 ± 0.10	$44.14_{\pm 0.13}$	$45.34_{\pm 0.15}$	$47.26_{\pm 0.03}$	47.98 ± 0.07	$46.56_{\pm 0.01}$	
		S0	31.95 ± 0.02	31.73+0.13	$32.17_{\pm 0.10}$	34.53 ± 0.02	33.53 ± 0.01	35.60 ± 0.05	
	VideoChat	S1	$45.10_{\pm 0.07}$	$46.24_{\pm 0.05}$	$44.01_{\pm 0.10}$	44.25 ± 0.09	$44.76_{\pm 0.02}$	$43.75_{\pm 0.16}$	
		S2	$45.53_{\pm 0.11}$	$42.90_{\pm 0.27}$	$48.49_{\pm 0.10}$	$45.57_{\pm 0.03}$	$47.20_{\pm 0.12}$	$44.05_{\pm 0.05}$	
		S0	$35.70_{\pm 0.06}$	$43.08_{\pm 0.00}$	$30.47_{\pm 0.09}$	$35.27_{\pm 0.01}$	$41.16_{\pm 0.00}$	$30.86_{\pm 0.01}$	
	VideoChat2	S1	$37.56_{\pm 0.07}$	$44.62_{\pm 0.00}$	$32.43_{\pm 0.10}$	$38.71_{\pm 0.10}$	$45.14_{\pm 0.13}$	$33.88_{\pm 0.08}$	
		S2	$49.07_{\pm 0.26}$	$54.72_{\pm 0.41}$	$44.47_{\pm 0.15}$	$48.86_{\pm 0.05}$	$57.12_{\pm 0.08}$	$42.68_{\pm 0.04}$	
		S0	$39.21_{\pm 0.14}$	$40.56_{\pm 0.15}$	$37.94_{\pm 0.12}$	$40.53_{\pm 0.33}$	$40.44_{\pm 0.24}$	$40.62_{\pm 0.43}$	
	mPLUG-Owl	S1	$45.80_{\pm 0.06}$	$47.49_{\pm 0.04}$	$44.22_{\pm 0.07}$	$47.97_{\pm 0.04}$	$49.33_{\pm 0.03}$	$46.69_{\pm 0.05}$	
		S2	$52.73_{\pm 0.13}$	$54.54_{\pm 0.13}$	$51.04_{\pm 0.13}$	$50.95_{\pm 0.06}$	$56.40_{\pm 0.11}$	$46.47_{\pm 0.18}$	
		S0	$40.71_{\pm 0.10}$	$41.38_{\pm 0.25}$	$40.07_{\pm 0.04}$	$43.45_{\pm 0.23}$	$43.24_{\pm 0.30}$	$43.66_{\pm 0.16}$	
	SALMONN	S1	$39.79_{\pm 0.03}$	$39.54_{\pm 0.01}$	$40.05_{\pm 0.06}$	$41.43_{\pm 0.13}$	$41.11_{\pm 0.03}$	$41.76_{\pm 0.22}$	
		S2	$47.96_{\pm 0.04}$	$50.20_{\pm 0.04}$	$45.92_{\pm 0.04}$	$48.24_{\pm 0.03}$	$52.24_{\pm 0.00}$	$44.82_{\pm 0.05}$	
		S0	$30.64_{\pm 0.06}$	$41.92_{\pm 0.00}$	$24.14_{\pm 0.08}$	$30.50_{\pm 0.05}$	$40.84_{\pm 0.13}$	$24.33_{\pm 0.03}$	
	Qwen-Audio	S1	$35.23_{\pm 0.10}$	$46.69_{\pm 0.15}$	$28.29_{\pm 0.08}$	$44.09_{\pm 0.00}$	$58.08_{\pm 0.00}$	$35.53_{\pm 0.00}$	
		S2	$38.13_{\pm 0.05}$	$49.42_{\pm 0.18}$	$31.04_{\pm 0.00}$	$41.14_{\pm 0.07}$	$53.71_{\pm 0.00}$	$33.34_{\pm 0.09}$	
		S0	$32.64_{\pm 0.03}$	$33.31_{\pm 0.01}$	$32.00_{\pm 0.05}$	$32.76_{\pm 0.03}$	$33.19_{\pm 0.06}$	$32.33_{\pm 0.00}$	
	Video-LLaVA	S1	$30.19_{\pm 0.02}$	$34.10_{\pm 0.03}$	$27.08_{\pm 0.05}$	$31.93_{\pm 0.11}$	$33.40_{\pm 0.19}$	$30.58_{\pm 0.04}$	
		S2	$47.07_{\pm 0.16}$	$48.58_{\pm 0.02}$	$45.66_{\pm 0.29}$	$49.21_{\pm 0.06}$	$53.95_{\pm 0.03}$	$45.23_{\pm 0.13}$	
		SO	$35.14_{\pm 0.14}$	$36.71_{\pm 0.15}$	$33.69_{\pm 0.14}$	$33.30_{\pm 0.04}$	$33.12_{\pm 0.06}$	$33.48_{\pm 0.03}$	
	LLaMA-VID	S1	$42.37_{\pm 0.03}$	$43.97_{\pm 0.04}$	$40.89_{\pm 0.03}$	$42.56_{\pm 0.08}$	$43.28_{\pm 0.11}$	$41.86_{\pm 0.04}$	
		S2	$51.25_{\pm 0.09}$	$52.71_{\pm 0.18}$	$49.87_{\pm 0.00}$	$52.01_{\pm 0.02}$	$57.30_{\pm 0.00}$	$47.61_{\pm 0.03}$	
		SO	$39.89_{\pm 0.18}$	$42.32_{\pm 0.21}$	$37.72_{\pm 0.15}$	$36.83_{\pm 0.30}$	$37.74_{\pm 0.27}$	$35.96_{\pm 0.33}$	
	Chat-UniVi	S1	$47.94_{\pm 0.19}$	$50.96_{\pm 0.20}$	$45.26_{\pm 0.18}$	$47.02_{\pm 0.00}$	$48.07_{\pm 0.00}$	$46.01_{\pm 0.00}$	
		S2	$53.08_{\pm 0.01}$	$53.68_{\pm 0.00}$	$52.50_{\pm 0.02}$	$53.86_{\pm 0.02}$	$58.54_{\pm 0.01}$	$49.86_{\pm 0.03}$	

1602 1603 1604

1605

O RELATIONSHIP BETWEEN DESCRIPTION LENGTH AND LABEL NUMBERS

This section further discusses the relationship between description length and the number of labels per sample, i.e., whether longer descriptions correlate with more labels. To this end, we compute their PCC scores. We observe that, for the human-only strategy, the PCC score is 0.3416, and for the human-LLM collaboration strategy, the PCC score is 0.2939. Therefore, although from the dataset level, the length of descriptions is related to the richness of labels (see Figure 5), these two metrics do not show a strong correlation at the sample level.

1612 1613

P COST OF GPT-BASED METRICS

1614

1615 This paper reports zero-shot performance, only focusing on the inference process. The cost of 1616 evaluating our OV-MERD dataset is about \$1 per evaluation, which may not seem high. However, 1617 for future work aimed at training frameworks to better address the OV-MER task, this cost will 1618 become prohibitive. For example, if we plan to train a model for 100 epochs, the evaluation cost 1619 will rise to \$1 × 100 epochs = \$100. If we intend to test N different parameter combinations and 1619 M different frameworks, the evaluation cost will further increase to \$100 × M × N. Moreover, we plan to expand the OV-MERD dataset in the future. This cost will further increase. Therefore, this paper explores alternatives to GPT-based metrics.

Q GPT-BASED VS. MATCHING-BASED METRICS

Table 13 provides raw scores for GPT- and matching-based metrics. See Section 5 for more analysis.

Table 13: GPT-based vs. matching-based metrics. " P_S ", " R_S ", " B_1 ", " B_4 ", "M', and " R_l " are abbreviations for Precision_s, Recall_s, BLEU₁, BLEU₄, METEOR, and ROUGE_l, respectively.

1630			English						Chinese									
1631	MLLM	MLLM L V A		GPT-based			Matching-based			GPT-based			Matching-based					
4000					Fs	P_S	Rs	B_1	\mathbf{B}_4	Μ	\mathbf{R}_l	Fs	Ps	Rs	B_1	\mathbf{B}_4	Μ	\mathbf{R}_l
1632	Qwen-Audio		×		38.13	49.42	31.04	21.87	06.55	21.65	20.81	41.14	53.71	33.34	27.64	12.07	26.09	25.24
1633	OneLLM		×		42.84	45.92	40.15	33.81	08.54	28.00	22.46	46.17	52.07	41.47	42.75	16.60	34.42	26.81
1634	Otter		′√	×	43.51	50.71	38.09	27.26	07.55	23.42	21.05	46.22	52.65	41.18	35.35	14.41	29.34	25.91
1007	Video-LLaMA		′√	×	44.73	44.14	45.34	28.76	06.41	31.22	20.41	47.26	47.98	46.56	34.88	12.13	37.61	24.25
1635	VideoChat		′√	×	45.53	42.90	48.49	26.44	05.41	30.58	19.11	45.57	47.20	44.05	31.36	10.86	37.48	22.57
1636	PandaGPT		′√		45.89	50.03	42.38	33.69	07.64	30.29	22.07	47.33	53.01	42.75	43.02	15.83	37.94	26.87
1637	Video-LLaVA		′√	×	47.07	48.58	45.66	33.48	08.25	29.68	22.34	49.21	53.95	45.23	42.72	15.97	36.87	26.90
1037	SALMONN		×		47.96	50.20	45.92	31.89	07.19	28.42	20.99	48.24	52.24	44.82	39.00	14.00	35.12	25.35
1638	VideoChat2		′√	×	49.07	54.72	44.47	31.60	08.10	26.61	21.65	48.86	57.12	42.68	41.18	16.15	33.54	26.80
1639	Video-ChatGPT		′√	×	50.52	54.03	47.44	32.64	07.65	30.25	22.01	54.73	61.15	49.52	41.96	15.50	38.18	26.35
1640	OneLLM		′√	×	50.52	55.93	46.06	32.19	08.10	28.44	22.25	51.44	56.43	47.26	41.31	15.15	35.15	25.98
1040	LLaMA-VID		′√	×	51.25	52.71	49.87	33.81	08.26	30.31	22.36	52.01	57.30	47.61	43.01	16.23	37.92	27.20
1641	mPLUG-Owl		′√	×	52.73	54.54	51.04	33.04	07.75	30.24	21.75	50.95	56.40	46.47	41.69	15.16	37.81	26.39
1642	Chat-UniVi			×	53.08	53.68	52.50	32.80	07.83	31.12	22.15	53.86	58.54	49.86	40.76	15.05	38.75	26.43
1040	GPT-4V		′√	×	55.51	48.52	64.86	39.40	18.41	43.67	32.60	57.21	54.61	60.07	45.45	29.08	53.76	40.37

In Table 13, we observe that there is no strong correlation between the GPT-based metrics and the matching-based metrics. To clarify this point, we use the following three sentences as examples:

1647 #1. The clue is "the weather is great". His emotion is "happy".

1648 #2. The clue is "the weather is bad". His emotion is "sad".

1650 #3. His emotion is "happy".

For matching-based metrics, we use $BLEU_1$ as an example. The $BLEU_1$ score between #1 and #2 is 0.8181, while the $BLEU_1$ score between #1 and #3 is 0.1738. Therefore, based on the $BLEU_1$ score, #1 is closer to #2. For LLM-based metrics, we first extract the emotion labels and compare their similarity, so #1 is closer to #3. This demonstrates that matching-based metrics are not suitable for evaluating emotion recognition performance.

1674 R EMOTION WHEEL

The emotion wheel provides psychologically based emotion grouping information. In this paper, we select five representative emotion wheels (W1 \sim W5) and use their grouping information for metric calculation. See Figure 21 for more details.

