

---

# EmCoBench: An Extensive Benchmark for General Emotion Comprehension

---

Yuxiang Lin<sup>1</sup> Jue Wang<sup>1,4</sup> Haomin Liang<sup>1</sup> Zebang Cheng<sup>1</sup>

Jun-Yan He<sup>2\*</sup> Zhi-Qi Cheng<sup>3\*</sup> Xiaojiang Peng<sup>1\*</sup> Alexander G. Hauptmann<sup>3</sup>

<sup>1</sup>Shenzhen Technology University <sup>2</sup>Alibaba Group <sup>3</sup>Carnegie Mellon University

<sup>4</sup>Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

<https://github.com/Lum1104/EmCoBench>

## Abstract

Affect computing is crucial in fields such as human-computer interaction, health-care, and market research, yet emotion’s ambiguity and subjectivity challenge current recognition techniques. We propose **Emotion Comprehension** (EC), a task that explains the reasons behind emotions, and create the **Emotion Comprehension Benchmark** (EmCoBench) using Vision Large Language Models (VLLMs) and a VLLM-assisted dataset construction method Coarse-to-Fine Self-Ask (CFSAs). EmCoBench includes 1,655 samples with 50 subsets and 78 emotion types. Experiments show limited proficiency of existing models in EC, with the best achieving 62.41% accuracy in the zero-shot setting and some performing lower than the text-only LLaMA-3 model (6.26%) in the caption-provided setting. Given its similarity to Visual Question Answering (VQA), EC can serve as a benchmark for assessing VLLMs’ fine-grained emotion understanding. Addressing EC’s challenges can lead to more empathetic AI systems, enhancing human-computer interaction and emotion-sensitive applications.

## 1 Introduction

Affect computing plays a crucial role across diverse domains [22], such as human-computer interaction (HCI) [20, 36, 41, 57], healthcare [12, 44, 46], and market research [5, 6, 45]. While recent work primarily focuses on emotion recognition by categorizing emotions into basic types, the inherent complexity and subjectivity of emotions make it challenging for individuals to accurately categorize their own emotions in complex situations. Despite this, these applications share a common goal: understanding the triggers of emotions. The variability of emotional experiences across individuals and contexts underscores the need to move beyond simple categorization and focus on the triggers and circumstances that lead to emotions, known as “Emotional Triggers”.

In response to this motivation, we propose the **Emotion Comprehension** (EC) task, which focuses on explaining emotional triggers rather than categorizing emotions. As shown in Figure 1, this task involves identifying the specific causes of emotional states for given individuals or scenes, aiming to provide valuable insights and solutions across various fields. Vision Large Language Models (VLLMs), known for their extensive world knowledge and explanatory abilities, are well-suited for emotion comprehension tasks [2, 8, 31, 34, 50, 33, 32, 26]. By framing the task in a Visual Question Answering (VQA) format, where VLLMs excel, we can better investigate the formation of emotions through visual and contextual data.

---

\*Co-corresponding authors: pengxiaojiang@sztu.edu.cn; junyanhe1989@gmail.com; zhiqic@cs.cmu.edu

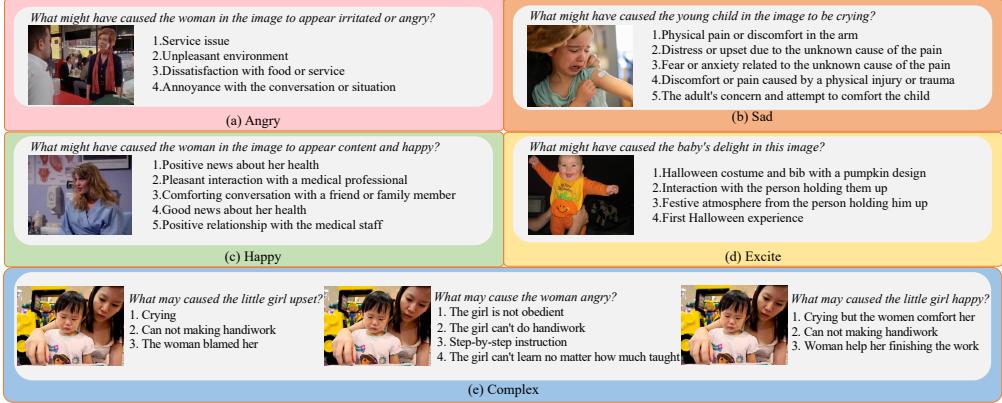


Figure 1: Figure (a-e) represent the different emotion categories under five scheme: angry, sad, excite, happy, and complex respectively, and the Emotion-trigger pairs.

To advance research in this area, we established the **Emotion Comprehension Benchmark** (Em-CoBench) (Figure 1), which includes 78 fine-grained emotions and 1,655 emotion comprehension samples, with 50 complex multifaceted emotions subset, such as combinations of happiness and sadness (Figure 1 (e)). We also proposed a VLLM-assisted data annotation scheme, the Coarse-to-Fine Self-Ask (CFSAs) method, which employs the Chain of Thought (CoT) approach [42, 37, 58, 3, 61, 60] to guide LLMs in preliminary annotation. Given the task’s similarity to VQA, this benchmark could serve as a future standard for evaluating the fine-grained emotion knowledge of VLLMs.

To thoroughly evaluate different models’ abilities in the EC task, we conducted a comprehensive assessment of commonly used open-source and closed-source models. We include 3 evaluating settings: 1) zero-shot to user questions, 2) combining image captions with user questions, and 3) reasoning with CoT. The experiments indicate that existing models still have significant shortcomings in interpreting emotions, with closed-source models generally outperforming open-source models. Interestingly, the evaluation results for basic and multifaceted emotions deviate from the expected pattern. The Claude-3 series, which performed best on the basic subset, underperforms the ChatGPT-4 series on the multifaceted complex subset. This phenomenon might be due to the Claude series’ tendency to provide definitive answers, whereas ChatGPT-4 demonstrates better capability in handling multifaceted reasoning.

Our contributions include: (1) proposing the Emotion Comprehension task, which focuses on identifying emotional triggers rather than merely classifying emotions, crucial for building more empathetic systems; (2) introducing the Coarse-to-Fine Self-Ask (CFSAs) method, a novel VLLM-assisted data annotation technique; (3) developing the EmCoBench dataset, which includes 78 fine-grained emotions and 1,655 emotion comprehension samples, with 50 multifaceted complex samples; and (4) extensively evaluating both open-source and closed-source models, highlighting the limitations of current models in interpreting emotions.

## 2 Related Work

### 2.1 Emotion Recognition

**Context Aware Emotion Recognition:** Context-aware emotion Recognition (CAER) goes beyond Facial Expression Recognition (FER), which focuses solely on perceiving emotion via the face [49, 48, 47, 62, 38, 29, 10]. CAER considers the emotional cues provided by background contexts, integrating facial and body language in a joint and boosting manner [23, 53, 51, 4, 43, 39, 28, 54, 59]. Various methods have contributed to this area. For example, [23] established the EMOTIC dataset and proposed a baseline that combines the body region and the whole image as the context. [24] proposed a dataset derived from movies with human social context. [53] built a context dictionary based on clusters of visual features to identify meaningful patterns using visual context. Additionally, [51] explored CAER with the commonsense knowledge from VLLMs, achieving state-of-the-art results on the EMOTIC [23] and CAER-S [24] datasets.

**Explainable Emotion Recognition with LLMs:** Large Language Models (LLMs) have been increasingly used for emotion recognition due to their extensive world knowledge [40, 35, 14, 25]. In NLP, [14] used Chain of Thought (CoT) prompting to recognize implicit emotions, while [25] designed a retrieval-based system for conversational emotion recognition. With the development of Vision LLMs (VLLMs) [34, 31, 33], researchers began exploring their use in emotion recognition [11]. However, due to a lack of emotion-related datasets, the zero-shot capability of LLMs in this area is limited. To address this, [52] employed visual instruction tuning for better VLLM performance, [51] used VLLMs to extract commonsense context, combining it with image data to train a transformer model, and [14] utilized the CoT approach to incrementally guide LLMs through emotion tasks.

Our work is not focus on Emotion Recognition, instead, we aim at interpreting the implicit emotion trigger behind human’s emotion, comprehending the formation of human’s emotion.

## 2.2 Humor Study

Humor is an integral part of human life and has been the focus of extensive research [7, 18, 17, 19, 9, 16, 55, 1, 15]. For instance, [7] examined the elements or characteristics within cartoon scenes that contribute to humor. Similarly, Memecap [18] compiled a dataset of 6.3K visual memes with visual metaphors to facilitate meme interpretation. [17] tested the ability of large language models (LLMs) to understand humor using a subset of the New Yorker Cartoon Caption Contest. Additionally, [19] introduced the Video Laugh Reasoning task to explain why people laugh in specific videos. [9] investigated pretrained LLMs’ ability to understand Chinese humor, which may also aid future research in humor generation. Building on these insights, our research takes a broader approach, moving beyond humor to understanding general emotions.

## 2.3 Chain of Thoughts Prompting

Chain of Thought (CoT) Prompting enhances problem-solving by breaking complex tasks into manageable, sequential steps, improving both accuracy and transparency [42, 37, 58, 3, 61, 60]. [42] proposed the Self-Ask method, where LLMs generate and answer their own sub-questions to solve a larger problem. [61] extended CoT to multimodal tasks, using a two-step approach where LLMs first generate a rationale and then use it for reasoning. [60] introduced a method for LLMs to solve tasks step-by-step in a one-by-one manner. We propose a Coarse-to-Fine Self-Ask method to guide VLLMs in assisting the EmCoBench annotation. This method progresses from general to scenario-specific perspectives, gradually deepening the understanding of emotional triggers.

## 3 Problem Definition

**Emotional Trigger:** We define an emotional trigger to any stimulus that provokes an emotional response in an individual (see Table 1). These triggers can include atmospheres, such as a dimly lit room or a bustling city street, and social interactions like conversations or arguments. Body movements and facial expressions also serve as triggers, where gestures or a smile can prompt specific emotions. Objects with sentimental value, performances in music or theater, outdoor activities like hiking, clothing choices, and sports events all play significant roles in evoking emotions.

While many of these triggers can be visually depicted, there are also implicit triggers that can evoke emotions beyond the image, such as an athlete’s adrenaline rush or the special feelings associated with people and scenes. We categorize these as “other” triggers. These varied triggers can deeply influence an individual’s emotional state, often subconsciously, shaping their reactions and interactions.

**Task Definition:** The Emotion Comprehension task aims to identify the implicit emotional trigger for a given input image  $x$  and emotion state  $e$  pair. These emotions can vary, being either positive (e.g., *happy*) or negative (e.g., *unhappy*) for the same individual. Mathematically, the Emotion Comprehension task  $\mathcal{F}$  can be formulated as follows:

$$\mathcal{F} : \mathcal{Q} \longrightarrow \mathcal{T} \quad (1)$$

Table 1: Emotional Trigger Types

Atmosphere	Social Interactions	Body Movements	Facial Expressions	Objects
Performances	Outdoor Activities	Clothing	Sports	Other

Table 2: This table illustrates how our CFSA Method comprehends a person’s excitement and joy during an LGBT event. **Pink** represent the generated caption, while **yellow** denotes the user query content, with both integrated. The **light orange** and trigger in the query are matched accordingly.

An Example of Comprehending Excitement		
User Question (VLLM)		
User Question Preprocessing (VLLM)	Question Parse	Who is the person user talks about? The man holding the rain bow flag. What is the user's demand? To know why the man is so <b>excited and joyful</b> .
	Question Reconstruction	The man holding the rainbow flag is a proud and enthusiastic participant in a parade, possibly celebrating LGBTQ+ pride or a similar event. He is waving the flag high in the air, showing his support and joy. <b>The user wants to know why the man is so excited and joyful.</b>
General Self-Ask Thinking (Question: Text-only LLM (GPT-4)) (Answer: VLLM)		1. How many people in this image, what are they doing? ... They are standing on top of a box ... are likely participating in a <b>parade or public event</b> . 2. What does “the person user mentioned” do in this image? Does he have interaction with others? ... is holding the rainbow flag and waving it high in the air... 3. What is the facial expression of the person mentioned? The facial expression of the person holding the flag is one of excitement and joy. 4. How about his body language? ... holding the flag is <b>expressive and energetic</b> , ... <b>active participation in the event</b> .
Scenario Specific Self-Ask Thinking (VLLM)		1. What specific ... the man holding the rainbow flag ... relates to his excitement and joy? The man is celebrating LGBTQ+ pride ... as indicated by the rainbow flag. ... participating in such an event is a way for him to express his <b>support for the LGBTQ+ community and celebrate their achievements and progress</b> . 2. What factors, such as personal experiences or beliefs, might contribute to the man’s enthusiastic participation in the parade or event? ... he may have <b>friends or family members who identify as LGBTQ+</b> ... 3. How does the man’s enthusiasm and joy impact the overall atmosphere and energy of the crowd and the event? The man’s enthusiasm ... <b>energy of the crowd and the event</b> ... 4. Are there any interactions or reactions from the crowd that might be influencing or being influenced by the man’s excitement and joy? ... provide him with a <b>sense of validation and encouragement</b> for his participation in the event.
Emotion Summarization (Emotional Triggers)		1. Celebrating LGBTQ+ pride or a similar event. 2. Supporting the LGBTQ+ community. 3. Expressing his personal beliefs and values. 4. Feeling a sense of unity and belonging with the crowd. 5. Being part of a positive and uplifting event. 6. Standing on top of a box.

where  $\mathcal{Q}$  is the image-emotion state pair set while  $\mathcal{T}$  denotes the set of emotional trigger. The goal of the task is to find the most probable triggers for the given input image-emotion state pair:

$$t = \mathcal{F}(q), \quad q \in \mathcal{Q}, \quad t \in \mathcal{T} \quad (2)$$

where image-emotion state pair  $q = (x, e)$ ,  $t$  is a text contain the meaning of emotional triggers.

This task shares similarities with “Explainable Multimodal Emotion Reasoning” (EMER) [30] as both aim to provide explanations. While EMER can support multi-classification, the nature of classification tasks limits its ability to output contradictory emotions.

## 4 Emotion Comprehension Benchmark

We refine the CAER-S [24] and EmoSet [56] context-aware emotion recognition datasets using the CFSA method, assisted with human annotation, to create a benchmark dataset. To the best of our knowledge, this is the first dataset focused on interpreting human emotions, featuring 1,615 basic EC and 50 complex EC samples with 78 fine-grained emotion types—surpassing other emotion-related datasets (Table 3).

### 4.1 VLLM-assisted Dataset Construction

#### 4.1.1 Coarse-to-Fine Self-Ask Annotation

We develop a Coarse-to-Fine Self-Ask (CFSA) method (Appendix Figure 4) to assist the emotion comprehension dataset annotating. CFSA involves breaking down complex, implicit user questions

into a series of simple VQA queries. Specifically, the VLLM assistance annotation process can be decoupled into four phases: 1) initial question preprocessing, 2) general self-ask thinking, 3) scenario self-ask thinking, and 4) emotion summarization. Notably, to ensure a high-quality EC dataset, four volunteers conducted a thorough manual review and detailed annotation of the entire dataset at all the phases of annotation.

**Initial Question Preprocessing.** To comprehensively and automatically capture the visual context of the image related to the emotion, we utilize a fixed and simple prompt to stimulate the LLM to complete and rich the prompt for visual questioning. Technically, we first parse the initial prompt as:  $s^{par} = \phi(s^{init})$ , where  $\phi$ ,  $s^{init}$ , and  $s^{par}$  represent the GPT-4, initial prompt, and the parsed question (prompt), respectively.  $s^{init}$  is constructed simply by a given emotion state  $e$  and prompt template. After the rich prompt  $s_{par}$  is produced, we collect the visual details rich reconstruction question as follows:

$$s_i^{rec} = llava(x_i, s^{par}) \quad x_i \in \mathcal{X} \quad (3)$$

where  $llava$  and  $x_i$  are the LLava-v1.6-34B (LLava-NEXT) [33] VLLM and the input image, respectively.  $s_i^{rec}$  denotes the reconstructed question produced by the LLava with given image  $x_i$ . The whole image dataset is denoted by  $\mathcal{X}$ .

Though VLLMs can describe images in detail, they often overlook emotional triggers due to limited emotion knowledge and reasoning. However, with the right question prompts, their strong VQA capabilities can help uncover these triggers. We introduce general self-ask and scenario self-ask methods to guide this process.

**General Self-Ask Thinking.** We let GPT-4 generate open-ended questions for the entire dataset. Afterward, we identified the four most frequently asked questions to prompt the VLLMs. Specifically, we compile these questions into a set  $\mathcal{S}^{gen} = \{s_1^{gen}, \dots, s_N^{gen}\}$  for all the images in our dataset. We then identify the four most frequently asked questions, denoted as  $\mathcal{S}^{freq} = \{s_1^{freq}, s_2^{freq}, s_3^{freq}, s_4^{freq}\}$ , and use these to prompt the VLLMs. And then the representative question set  $\mathcal{S}^{freq}$ . The  $\mathcal{S}^{freq}$  is mainly focused on 4 aspects:

- **Number of people in the image:** Provides context for the individuals' emotional states, as those around them may influence their emotions.
- **Activity and Interactions:** Understanding individuals' actions and interactions with others can reveal more about their emotional states and the scene's context.
- **Facial Expressions:** Key indicators of emotions, providing insight into human's feelings.
- **Body Language:** Conveys mood, intentions, and complementing facial expressions.

These four types of questions are further leveraged to query the visual details,  $a_i^{gen} = llava(x_i, s_i^{freq})$ , where  $a_i^{gen}$  is the answer provided by  $llava$ . The  $a_i^{gen}$  is further collected into an answer set  $\mathcal{A}^{gen} = \{a_1^{gen}, a_2^{gen}, a_3^{gen}, a_4^{gen}\}$ .

**Scenario Self-Ask Thinking.** Going a step further, we provide the VLLM model with the image example, user question  $s^{query}$ , reconstructed question  $s^{rec}$ , and the general self-ask question-answer pairs  $\mathcal{S}^{freq}$ ,  $\mathcal{A}^{gen}$  to produce the rich scenario details description,  $\mathcal{S}_i^{sce} = llava(x_i, [s^{query}, s^{rec}, \mathcal{S}^{freq}, \mathcal{A}^{gen}])$ , where  $[.]$  denotes the concatenate operation, and  $\mathcal{S}_i^{sce}$  denotes the scenario self-ask question set. Following this, the scenario self-ask answer set  $\mathcal{A}^{sce} = \{a_1^{sce}, a_2^{sce}, a_3^{sce}, a_4^{sce}\}$  is generated by  $a_i^{sce} = llava(x_i, [s^{query}, s^{rec}, \mathcal{S}^{freq}, \mathcal{A}^{gen}, \mathcal{S}_i^{sce}])$ , where  $a_i^{sce}$  is the scenario self-ask answer.

**Emotion Summarization.** After the general and scenario self-ask thinking, the critical factors impacting human emotion are comprehensively investigated, therefore, the emotion triggers can be summarized by the LLM model easily. To economically and practically summarize the emotional triggers, we leverage the recent powerful open-source LLM model, LLaMA-3, to extract them with an in-context learning scheme from all the LLava outputs.

**Human in the Loop Annotation** We use CDSA as a baseline annotation and employ LLaMA-3 for emotional trigger extraction. An example of the annotation process is depicted in Table 2. Following this, four volunteers conducted a thorough manual review and detailed annotation of the entire dataset at all the phases of annotation, with three core goals: 1) Remove hallucinations generated by the

VLLMs (Appendix B.1), 2) Add more commonsense knowledge to the emotion comprehension process (Appendix B.2), and 3) Curate the dataset by removing unnecessary emotional triggers.

## 4.2 Evaluation Metric

**Emotional Trigger Recall and Long-term Coherence** Given the subjective nature of emotions, multiple triggers could elicit a particular response, and some may be missed despite thorough reviews. Therefore, we use *Recall* as one evaluation metric. If the model’s understanding overlaps with our ground truth, it is considered correct. An emotional trigger identified by the model is a true positive if it overlaps with any part of our ground truth annotations; otherwise, it is a false negative. Additionally, *Long-term Coherence* in the context of emotion comprehension evaluates a model’s ability to maintain consistent emotional and thematic understanding throughout extended text. This metric is crucial for tasks where the emotional narrative or flow must remain logical and coherent over multiple sentences or paragraphs. We implemented this using LLaMA-3 and ChatGPT-3.5 (gpt-3.5-turbo-0125) for recall metric calculations and BERT [13] embeddings for the long-term coherence.

## 4.3 Dataset Overview

We chose context-aware emotion recognition datasets with rich background information, as emotion comprehension aims to uncover emotional triggers. Facial expression recognition datasets, including close-up facial data, are unsuitable for this purpose. The CAER-S dataset [24], derived from movie clips, includes a variety of life scenarios portraying seven emotions: *angry, disgust, fear, happy, neutral, sad, and surprise*. EmoSet [56], sourced from internet searches, is annotated with both positive emotions (*amusement, awe, contentment, excitement*) and negative emotions (*anger, disgust, fear, sadness*). Considering the inherent uncertainty of emotions and the need for manual review and filtering, we selected four emotions for our initial attempt to construct an emotion comprehension dataset. These include *happy* and *angry* from CAER-S, and *excitement* and *sadness* from EmoSet.

Table 3: Comparison of Various Emotion Datasets. The table highlights the differences in datasets used for emotion-related tasks. ER stands for Emotion Recognition, EMER stands for Explainable Multimodal Emotion Recognition Reasoning, and EC stands for Emotion Comprehension.

Dataset	Task	Annotator	Emotion Types	Explainable	Has Complex Label
CAER-S [24]	ER	6	7	✗	✗
DFEW [21]	ER	3	7	✗	✗
RAF-DB [27]	ER	315	7	✗	✗
HECO [54]	ER	13	8	✗	✗
EMOTIC [23]	ER	-	26	✗	✗
EmoSet [56]	ER	10	8	✓	✗
MER2023(EMER) [30]	EMER	6	7	✓	✗
Ours	EC	4 (VLLM-assisted)	78	✓	✓

## 4.4 Data Analysis

Our benchmark extensively explores human emotions, categorizing them into four primary groups: angry, sad, excited, and happy. Each primary category is divided into fine-grained emotions, totaling 78. The statistics on these fine-grained emotions within each primary category are illustrated in Table 4. For instance, the anger category includes emotions such as *annoyed, agitated*, and *upset*, capturing varying intensities of anger. The sadness category includes emotions like *forlorn* and *contemplative*, highlighting different depths of sadness. The excitement category features emotions such as *delight* and *thrill*, reflecting different degrees of enthusiasm. The happiness category includes emotions like *lighthearted*, offering insights into various states of joy. The multifaceted complex subset consists of 50 samples, each interpreted from more than two perspectives.

**Comparison.** Table 3 compares our dataset with other emotion-related datasets. Our dataset is notable for its 78 fine-grained emotions and its interpretability. The complex subset features intricate emotions, including difficult labels absent in other datasets. Additional visualizations of our complex EC subset can be found in Appendix A.2.

**Emotion Trigger Distribution** The emotion triggers are divided into 10 categories, with definitions provided in Table 1. We also present the distribution of emotion triggers across these

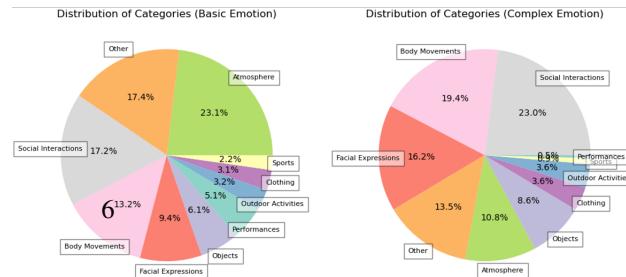


Figure 2: Visualization of the numbers of emotional triggers

categories, depicted in Figure 4.4. "Atmosphere" and "Others" are the top two triggers in basic emotions, while "Social Interaction" and "Body Movements" top the list in complex emotions.

Table 4: Fine-grained emotional breakdown within primary emotional categories.

	Primary	Fine-grained
Negative	Angry	angry (191), annoyed (85), agitated (64), upset (69), irritated (60), outraged (37), infuriated (22), hostile (24), concerned (11), frustrated (11), serious (4), displeased (6), mad (1), surprised (3), shocked (2), exhibit (1)
	Sad	sadness (241), forlorn (82), contemplative (71), unhappy (45), disheartened (27), dismal (28), solemn (23), sorrowful (20), somber (11), distress (13), miserable (13), discontent (5), upset (4), disappointment (3), distraught (2), displeased (2), frown (1), weary (1), frustration (1), loneliness (1), tragic, disappointed (1), melancholic (1), pain (1), injury (1)
Positive	Excite	excited (152), delight (63), thrill (59), inspired (68), stimulate (49), incite (42), spur (38), smile (45), happy (18), raised (23), joyful (29), fascinating (35), enjoying (8), brightly (6), spark (8), enthusiasm (6), funny (3), intense (2), pleasant (2), feathery (1)
	Happy	happy (61), smile (60), lighthearted (29), radiant (27), contented (46), pleased (22), spirited (21), cheerful (18), exhilarated (13), glad (10), blissful (12), energetic (15), joyful (5), optimistic (4), enjoying (3), positive (2), surprised (2)

## 5 Experiments

In this section, we evaluate the performance of both prominent open-source models and a proprietary API on our benchmark. We employed four different modes to assess the models' capabilities in emotion comprehension. Following the evaluation, we present a comprehensive analysis of how each model performs in terms of emotion comprehension. In the experiment, we adopted three modes to implement the emotion comprehension algorithms of various models. 1) *User Question*, involves zero-shot testing of the model using the user questions, evaluating the abilities of how models deal with human questions. 2) *User Question and Caption*, involves inputting questions that are parsed and restructured based on Section 4 to make the user's query more specific. Additionally, we tested the text-only model LLaMA-3 with the provided caption. 3) *CFSAs setting*, uses the CFSAs implemented by the LLaVA-NEXT 34B model to provide responses, showing the gap between human performance and the vllm assisted annotation dataset. And 4) *User Question and CoT setting*, inspired by the Chain of Thought (CoT) approach, involves adding the prompt "let's think step by step" after the user's question to encourage the model to reason through the problem methodically.

### 5.1 Model Performance Overview

The overall performance comparison of LLMs across basic emotion comprehension (Table 5) and complex emotion comprehension (Table 6), reveals the distinct strengths and weaknesses of each model. In the open-source models, the LLaVA series and MiniGPT-v2 comprehend emotions well, while Qwen-VL-Chat consistently attains the lowest scores. Video-LLaVA and Otter perform moderately, although Otter shows a notable weakness in handling the excitement category. Overall, closed-source models such as ChatGPT-4 and the Claude-3 series outperform open-source models when processing user questions alone, though the Qwen-vl-plus model's performance remains subpar. Notably, the closed-source Claude-3 series (claude-3-sonnet-20240229, claude-3-haiku-20240307) excels in the basic emotion comprehension setting, securing the highest overall scores and demonstrating strong emotion comprehension abilities (Table 5 Close-source). However, in the complex setting, its scores are lower than those of the ChatGPT-4 series (Table 6). Table 7 reports the long-term coherence scores between models, showing that their scores are close, demonstrating a consistent ability to maintain context over comprehending emotion.

The CFSAs scores reveal that while emotion comprehension is challenging for models, they can successfully identify 68% of emotional triggers. By converting the problem into a series of simple VQA tasks, the scores show a significant improvement. Additionally, the scores indicate that the

Table 5: Basic Emotion comprehension performance of Open-Source/Close-Source Language Models, with evaluation scores presented for each subclass according to the LLaMA-3/ChatGPT criteria.

Models	Happy	Angry	Sadness	Excitement	Overall
<b>User Question</b>					
Qwen-VL-Chat	32.09/39.68	22.32/26.10	30.64/33.88	25.02/36.32	26.45/33.65
Video-LLaVA	55.55/53.28	40.42/36.97	50.62/45.25	51.78/52.23	49.26/47.06
MiniGPT-v2	52.78/51.80	<b>47.10/47.76</b>	<b>60.47/58.14</b>	50.78/53.66	52.89/53.59
Otter	45.63/49.25	42.53/43.07	47.67/46.19	39.47/48.30	42.81/46.64
LLaVA-1.5 (13B)	<b>59.01/57.52</b>	45.44/41.88	55.16/48.64	<b>57.46/58.73</b>	<b>54.37/52.20</b>
LLaVA-NEXT (7B)	54.16/49.24	43.71/39.87	53.29/46.52	58.90/53.06	53.82/48.18
LLaVA-NEXT (13B)	57.17/55.18	43.16/37.93	54.16/45.42	59.38/55.29	54.33/48.79
LLaVA-NEXT (34B)	54.50/51.03	38.96/35.65	51.10/47.21	51.77/52.04	49.03/47.13
<b>User Question &amp; Caption</b>					
Qwen-VL-Chat	41.94/46.34	32.71/31.91	41.82/44.16	38.65/43.84	38.47/41.54
Video-LLaVA	56.77/58.79	43.65/43.86	54.25/55.12	55.35/59.42	52.63/54.85
MiniGPT-v2	55.11/60.04	47.95/51.00	<b>62.29/64.24</b>	51.55/57.90	54.05/58.37
Otter	48.97/54.67	34.22/37.12	34.57/37.55	35.27/42.99	35.62/40.85
LLaVA-1.5 (13B)	57.91/58.46	43.75/40.72	55.47/51.46	56.42/59.42	53.55/53.13
LLaVA-NEXT (7B)	<b>64.32/61.00</b>	48.60/46.74	58.75/53.00	<b>62.99/59.39</b>	58.80/54.97
LLaVA-NEXT (13B)	61.99/61.95	<b>48.84/46.85</b>	59.62/55.18	62.17/59.95	<b>58.60/55.92</b>
LLaVA-NEXT (34B)	57.51/62.73	46.47/47.87	58.35/55.84	60.17/59.64	56.60/56.24
LLaMA-3 (8B) ( <b>Text Only</b> )	52.36/50.73	34.78/32.71	52.29/46.87	43.62/42.06	44.73/41.94
<b>User Question &amp; CoT</b>					
Qwen-VL-Chat	41.99/44.46	34.62/31.06	43.64/39.30	32.78/40.04	36.79/38.18
Video-LLaVA	51.42/47.63	42.68/35.65	56.77/46.29	53.01/46.98	51.81/44.42
MiniGPT-v2	56.36/57.58	47.71/48.32	<b>59.46/56.79</b>	50.21/52.39	52.67/53.08
Otter	49.97/51.91	43.23/43.71	50.15/46.86	42.30/47.16	45.17/46.61
LLaVA-1.5 (13B)	<b>59.12/56.94</b>	40.97/34.44	53.07/45.66	54.16/54.36	51.34/47.80
LLaVA-NEXT (7B)	54.74/52.04	44.61/41.93	52.69/47.63	52.78/47.60	51.14/46.66
LLaVA-NEXT (13B)	50.91/50.35	42.21/38.81	54.66/49.42	51.64/49.39	50.47/47.21
LLaVA-NEXT (34B)	52.17/49.55	<b>48.35/44.45</b>	55.97/50.55	<b>55.29/53.46</b>	<b>53.84/50.50</b>
CFSAs (LLaVA-NEXT (34B))	69.68/68.72	61.08/61.14	68.39/69.46	72.63/70.31	68.81/68.04
<b>Close-source Models</b>					
Qwen-vl-plus	29.05/27.22	23.58/17.89	38.35/30.08	30.09/26.87	31.00/25.90
ChatGPT-4V	52.30/55.74	48.93/48.57	45.00/44.42	46.38/49.90	46.86/48.58
ChatGPT-4o	52.94/50.78	42.12/35.33	49.79/46.42	53.48/54.53	49.99/47.93
Claude-3-haiku	<b>59.20/60.28</b>	<b>49.87/49.84</b>	<b>67.21/63.26</b>	<b>67.55/68.10</b>	<b>63.24/62.41</b>
Claude-3-sonnet	44.58/44.45	38.95/42.86	55.98/54.40	61.41/62.24	54.10/54.89

VLLM assisted annotation, still lag behind human-level annotations, highlighting the considerable effort by our manually labeling.

## 5.2 Abilities Comparison

In the direct User Question setting, all models scored relatively low. After adding the Caption, the scores of all models improved, but it is notable that the Otter model’s overall score decreased by approximately 7%. MiniGPT-v2 scored higher in the Angry and Sadness categories, while the LLaVA series models performed better overall, with LLaVA v1.5 (13B) achieving the highest scores, particularly excelling in the Happy and Excitement categories. Interestingly, increasing the model size did not lead to better performance; the 34B model’s scores even declined. The Qwen-VL-Chat model performed poorly across all emotion categories, and its performance with Caption and user questions was even worse than that of the text-only LLaMA-3 model.

We introduced “let’s think step by step” as part of the input, leveraging the CoT approach to improve model performance. Results show that this method consistently outperformed the direct User Question setting, indicating the complexity of the emotion comprehension task. Detailed reasoning and step-by-step responses proved more effective than direct answers, helping uncover more emotional triggers. This finding aligns with our observations using the CFSAs method, where models better identified emotion triggers through detailed, step-by-step analysis.

Table 6 shows the performance of models on the multifaceted complex subset. The performance of open-source models is similar to that on the basic subset, but their scores are significantly lower. This subset evaluates the models’ abilities in multifaceted emotional reasoning and empathy. Even the highest-scoring model, LLaVA-1.5, only achieves 38.10/39.53 points, while is notably close to the best-performing closed-source model, ChatGPT-4. Notably, the Claude-3 series, which performed best on the basic subset, does not achieve state-of-the-art results on this subset and even scores lower

Table 6: Evaluation of the complex emotion comprehension ability among the VLLMs.

Models	Recall
<b><i>Open-Souce</i></b>	
Qwen-VL-Chat	22.00/32.40
Video-LLaVA	30.90/32.27
MiniGPT-v2	35.10/36.00
Otter	27.90/33.23
LLaVA-1.5 (13B)	<u>38.10/39.53</u>
LLaVA-NEXT (7B)	38.71/33.50
LLaVA-NEXT (13B)	39.16/33.60
LLaVA-NEXT (34B)	35.37/33.10
<b><i>Close-Source</i></b>	
Qwen-vl-plus	20.37/19.60
Claude-3-haiku	24.00/24.77
Claude-3-sonnet	21.37/22.40
ChatGPT-4V	28.00/30.60
ChatGPT-4o	<b>39.27/39.57</b>

Table 7: Metric of Long-term Coherence between

VLLMs in user question setting.

Models	Coherence
<b><i>Open-Souce</i></b>	
Qwen-VL-Chat	84.49
Video-LLaVA	84.89
MiniGPT-v2	84.70
Otter	<u>85.03</u>
LLaVA-1.5 (13B)	84.50
LLaVA-NEXT (7B)	81.02
LLaVA-NEXT (13B)	81.09
LLaVA-NEXT (34B)	84.96
<b><i>Close-Source</i></b>	
Qwen-vl-plus	83.00
Claude-3-haiku	<b>85.98</b>
Claude-3-sonnet	84.53
ChatGPT-4V	81.97
ChatGPT-4o	80.65

than some open-source models. This indicates that while the Claude-3 series excels at basic emotion comprehension, its ability to handle more complex, multifaceted emotional reasoning is less effective.

## 6 Conclusion and Discussion

In this paper, we introduce the Emotion Comprehension (EC) task, which focuses on identifying emotional triggers rather than merely labeling emotions. We establish the Emotion Comprehension Benchmark (EmCoBench) dataset using a VLLM-assisted construction method, Coarse-to-Fine Self-Ask (CFSAs), consisting of 1,655 samples, including a well-annotated multifaceted complex subset of 50 samples. Extensive experiments evaluating commonly used open-source and closed-source models demonstrate that these models have limited proficiency in the EC task. Some are even performing lower than the text-only LLaMA-3 model in the caption-provided setting.

Emotion comprehension is an open-world problem due to the unique personal experiences of each individual. We encourage the community to explore new approaches for this task to achieve a more comprehensive understanding of emotions. This task not only enhances emotion comprehension but also provides a novel metric for evaluating the emotional intelligence of VLLMs. By considering emotions from multiple perspectives, it aids in analyzing implicit emotions and advances the field of emotion recognition. Additionally, our annotation method facilitates emotion reasoning dataset development and provides a valuable resource for multi-turn dialogue in emotion research.

## Acknowledgments and Disclosure of Funding

We would like to thank...

## References

- [1] Issa Annamoradnejad and Gohar Zoghi. Colbert: Using bert sentence embedding for humor detection. *arXiv preprint arXiv:2004.12765*, 1(3), 2020.
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [3] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczek, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- [4] Uttaran Bhattacharya, Trisha Mittal, Rohan Chandra, Tanmay Randhavane, Aniket Bera, and Dinesh Manocha. Step: Spatial temporal graph convolutional networks for emotion perception from gaits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 1342–1350, 2020.

- [5] Erik Cambria, Dipankar Das, Sivaji Bandyopadhyay, and Antonio Feraco. Affective computing and sentiment analysis. *A practical guide to sentiment analysis*, pages 1–10, 2017.
- [6] Delphine Caruelle, Poja Shams, Anders Gustafsson, and Line Lervik-Olsen. Affective computing in marketing: practical implications and research opportunities afforded by emotionally intelligent machines. *Marketing Letters*, 33(1):163–169, 2022.
- [7] Arjun Chandrasekaran, Ashwin K Vijayakumar, Stanislaw Antol, Mohit Bansal, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. We are humor beings: Understanding and predicting visual humor. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4603–4612, 2016.
- [8] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.
- [9] Yuyan Chen, Zhixu Li, Jiaqing Liang, Yanghua Xiao, Bang Liu, and Yunwen Chen. Can pre-trained language models understand chinese humor? In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pages 465–480, 2023.
- [10] Zebang Cheng, Yuxiang Lin, Zhaoru Chen, Xiang Li, Shuyi Mao, Fan Zhang, Daijun Ding, Bowen Zhang, and Xiaojiang Peng. Semi-supervised multimodal emotion recognition with expression mae. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 9436–9440, 2023.
- [11] Zebang Cheng, Fuqiang Niu, Yuxiang Lin, Zhi-Qi Cheng, Bowen Zhang, and Xiaojiang Peng. Mips at semeval-2024 task 3: Multimodal emotion-cause pair extraction in conversations with multimodal language models. *arXiv preprint arXiv:2404.00511*, 2024.
- [12] Ronald E Dahl and Allison G Harvey. Sleep in children and adolescents with behavioral and emotional disorders. *Sleep medicine clinics*, 2(3):501–511, 2007.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [14] Hao Fei, Bobo Li, Qian Liu, Lidong Bing, Fei Li, and Tat-Seng Chua. Reasoning implicit sentiment with chain-of-thought prompting. *arXiv preprint arXiv:2305.11255*, 2023.
- [15] Md Kamrul Hasan, Sangwu Lee, Wasifur Rahman, Amir Zadeh, Rada Mihalcea, Louis-Philippe Morency, and Ehsan Hoque. Humor knowledge enriched transformer for understanding multimodal humor. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 12972–12980, 2021.
- [16] Md Kamrul Hasan, Wasifur Rahman, Amir Zadeh, Jianyuan Zhong, Md Iftekhar Tanveer, Louis-Philippe Morency, et al. Ur-funny: A multimodal language dataset for understanding humor. *arXiv preprint arXiv:1904.06618*, 2019.
- [17] Jack Hessel, Ana Marasović, Jena D Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff, and Yejin Choi. Do androids laugh at electric sheep? humor “understanding” benchmarks from the new yorker caption contest. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 688–714, 2023.
- [18] EunJeong Hwang and Vered Shwartz. Memecap: A dataset for captioning and interpreting memes. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1433–1445, 2023.
- [19] Lee Hyun, Kim Sung-Bin, Seungju Han, Youngjae Yu, and Tae-Hyun Oh. Smile: Multi-modal dataset for understanding laughter in video with language models. *arXiv preprint arXiv:2312.09818*, 2023.

- [20] Shilpi Jain, Sriparna Basu, Arghya Ray, and Ronnie Das. Impact of irritation and negative emotions on the performance of voice assistants: Netting dissatisfied customers' perspectives. *International Journal of Information Management*, 72:102662, 2023.
- [21] Xingxun Jiang, Yuan Zong, Wenming Zheng, Chuangao Tang, Wanchuang Xia, Cheng Lu, and Jiateng Liu. Dfew: A large-scale database for recognizing dynamic facial expressions in the wild. In *Proceedings of the 28th ACM international conference on multimedia*, pages 2881–2889, 2020.
- [22] Smith K. Khare, Victoria Blanes-Vidal, Esmaeil S. Nadimi, and U. Rajendra Acharya. Emotion recognition and artificial intelligence: A systematic review (2014–2023) and research recommendations. *Information Fusion*, 102:102019, 2024.
- [23] Ronak Kosti, Jose M Alvarez, Adria Recasens, and Agata Lapedriza. Emotic: Emotions in context dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 61–69, 2017.
- [24] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and Kwanghoon Sohn. Context-aware emotion recognition networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10143–10152, 2019.
- [25] Shanglin Lei, Guanting Dong, Xiaoping Wang, Keheng Wang, and Sirui Wang. Instructerc: Reforming emotion recognition in conversation with a retrieval multi-task llms framework. *arXiv preprint arXiv:2309.11911*, 2023.
- [26] 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*, 2023.
- [27] Shan Li and Weihong Deng. Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. *IEEE Transactions on Image Processing*, 28(1):356–370, 2019.
- [28] Weixin Li, Xuan Dong, and Yunhong Wang. Human emotion recognition with relational region-level analysis. *IEEE Transactions on Affective Computing*, 14(1):650–663, 2021.
- [29] Yande Li, Mingjie Wang, Minglun Gong, Yonggang Lu, and Li Liu. Fer-former: Multi-modal transformer for facial expression recognition. *arXiv preprint arXiv:2303.12997*, 2023.
- [30] Zheng Lian, Licai Sun, Mingyu Xu, Haiyang Sun, Ke Xu, Zhuofan Wen, Shun Chen, Bin Liu, and Jianhua Tao. Explainable multimodal emotion reasoning. *arXiv preprint arXiv:2306.15401*, 2023.
- [31] 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.
- [32] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- [33] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024.
- [34] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [35] Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. Generated knowledge prompting for commonsense reasoning. *arXiv preprint arXiv:2110.08387*, 2021.
- [36] Yong Ma, Heiko Drewes, and Andreas Butz. How should voice assistants deal with users' emotions? *arXiv preprint arXiv:2204.02212*, 2022.
- [37] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

- [38] Jiawei Mao, Rui Xu, Xuesong Yin, Yuanqi Chang, Binling Nie, and Aibin Huang. Poster++: A simpler and stronger facial expression recognition network. *arXiv preprint arXiv:2301.12149*, 2023.
- [39] Trisha Mittal, Pooja Guhan, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. Emoticon: Context-aware multimodal emotion recognition using frege’s principle. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14234–14243, 2020.
- [40] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [41] Emmi Parviainen and Marie Louise Juul Søndergaard. Experiential qualities of whispering with voice assistants. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI ’20, page 1–13, New York, NY, USA, 2020. Association for Computing Machinery.
- [42] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*, 2022.
- [43] Shulan Ruan, Kun Zhang, Yijun Wang, Hanqing Tao, Weidong He, Guangyi Lv, and Enhong Chen. Context-aware generation-based net for multi-label visual emotion recognition. In *2020 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE Computer Society, 2020.
- [44] Carolyn Saarni, Joseph J Campos, Linda A Camras, and David Witherington. Emotional development: Action, communication, and understanding. *Handbook of child psychology*, 3, 2007.
- [45] Gautam Srivastava and Surajit Bag. Modern-day marketing concepts based on face recognition and neuro-marketing: a review and future research directions. *Benchmarking: An International Journal*, 31(2):410–438, 2024.
- [46] Edward Z Tronick. Emotions and emotional communication in infants. *Parent-infant psychodynamics*, pages 35–53, 2018.
- [47] Thanh-Hung Vo, Guee-Sang Lee, Hyung-Jeong Yang, and Soo-Hyung Kim. Pyramid with super resolution for in-the-wild facial expression recognition. *IEEE Access*, 8:131988–132001, 2020.
- [48] Kai Wang, Xiaojiang Peng, Jianfei Yang, Shijian Lu, and Yu Qiao. Suppressing uncertainties for large-scale facial expression recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6897–6906, 2020.
- [49] Kai Wang, Xiaojiang Peng, Jianfei Yang, Debin Meng, and Yu Qiao. Region attention networks for pose and occlusion robust facial expression recognition. *IEEE Transactions on Image Processing*, 29:4057–4069, 2020.
- [50] Xinshun Wang, Zhongbin Fang, Xia Li, Xiangtai Li, Chen Chen, and Mengyuan Liu. Skeleton-in-context: Unified skeleton sequence modeling with in-context learning. *arXiv preprint arXiv:2312.03703*, 2023.
- [51] Alexandros Xenos, Niki Maria Foteinopoulou, Ioanna Ntinou, Ioannis Patras, and Georgios Tzimiropoulos. Vllms provide better context for emotion understanding through common sense reasoning. *arXiv preprint arXiv:2404.07078*, 2024.
- [52] Hongxia Xie, Chu-Jun Peng, Yu-Wen Tseng, Hung-Jen Chen, Chan-Feng Hsu, Hong-Han Shuai, and Wen-Huang Cheng. Emovit: Revolutionizing emotion insights with visual instruction tuning. *arXiv preprint arXiv:2404.16670*, 2024.

- [53] Dingkang Yang, Zhaoyu Chen, Yuzheng Wang, Shunli Wang, Mingcheng Li, Siao Liu, Xiao Zhao, Shuai Huang, Zhiyan Dong, Peng Zhai, et al. Context de-confounded emotion recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19005–19015, 2023.
- [54] Dingkang Yang, Shuai Huang, Shunli Wang, Yang Liu, Peng Zhai, Liuzhen Su, Mingcheng Li, and Lihua Zhang. Emotion recognition for multiple context awareness. In *European Conference on Computer Vision*, pages 144–162. Springer, 2022.
- [55] Dify Yang, Alon Lavie, Chris Dyer, and Eduard Hovy. Humor recognition and humor anchor extraction. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 2367–2376, 2015.
- [56] Jingyuan Yang, Qirui Huang, Tingting Ding, Dani Lischinski, Danny Cohen-Or, and Hui Huang. Emoset: A large-scale visual emotion dataset with rich attributes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20383–20394, 2023.
- [57] Xi Yang, Marco Aurisicchio, and Weston Baxter. Understanding affective experiences with conversational agents. In *proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–12, 2019.
- [58] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [59] Minghui Zhang, Yumeng Liang, and Huadong Ma. Context-aware affective graph reasoning for emotion recognition. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 151–156. IEEE, 2019.
- [60] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- [61] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.
- [62] Ce Zheng, Matias Mendieta, and Chen Chen. Poster: A pyramid cross-fusion transformer network for facial expression recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3146–3155, 2023.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
  - (b) Did you describe the limitations of your work? **[Yes]** See Section 6.
  - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]** Our work is focus on general emotion comprehension, which may not have the negative societal impacts.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** Our paper does not contain any theoretical derivations.
  - (b) Did you include complete proofs of all theoretical results? **[N/A]** Our paper does not contain any theoretical derivations.
3. If you ran experiments (e.g. for benchmarks)...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]**
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]**
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[N/A]**
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[N/A]** We provide various baseline VLLMs method that require different computing resources.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Section 4.3
  - (b) Did you mention the license of the assets? **[Yes]**
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[N/A]**
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]**
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]** We did not use crowdsourcing.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

## A Emotion Comprehension Dataset

### A.1 Basic Emotion Comprehension Dataset



Figure 3: Visualization of fine-grained emotions in user questions from our dataset.

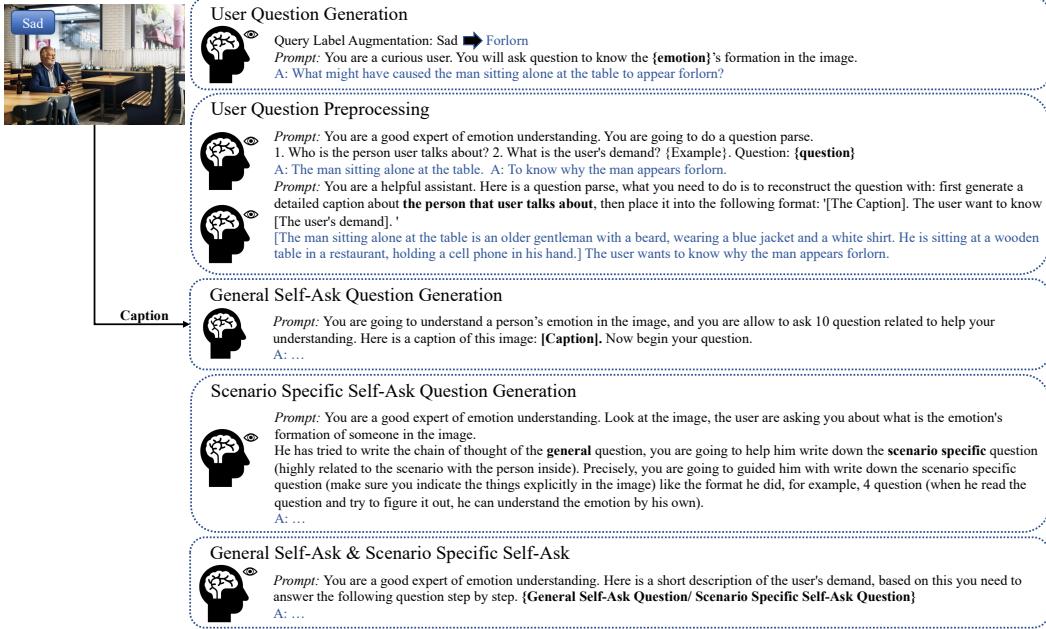


Figure 4: Pipeline of the VLLM-assisted dataset construction.

**Table 8: Visualization of basic emotion comprehension dataset, an image is corresponded to one user questions.**

**Examples of the Basic Emotion Comprehension Dataset**

User Question Emotional Trigger	
User Question Emotional Trigger	<p><i>What led to the formation of the arouse to the man in this image?</i></p> <p>1. Climbing a steep, snow-covered slope. 2. Physical effort and concentration. 3. Potential hazards and challenges. 4. Cold environment. 5. Determination to reach the goal.</p>
User Question Emotional Trigger	<p><i>What do you think might have caused the person's delight as they look out the window?</i></p> <p>1. Snowy scene outside the car. 2. Smile on her face. 3. Enjoying the view. 4. Serenity of the winter environment. 5. Excitement of experiencing a snowy day. 6. Personal or emotional connections to snowy weather or winter scenes. 7. Fresh snowfall, brightness of the snow reflecting sunlight, or peacefulness of the scene.</p>
User Question Emotional Trigger	<p><i>What do you think might have caused the man holding the box in the image to become lighthearted?</i></p> <p>1. Holding the "Überweiss" box. 2. Smiling. 3. Friendly and approachable body language. 4. Positive and relaxed atmosphere of the laundry room. 5. Interaction with others in the laundry room.</p>
User Question Emotional Trigger	<p><i>What might have caused the woman in the image to appear content and happy?</i></p> <p>1. Positive news about her health. 2. Pleasant interaction with a medical professional. 3. Comforting conversation with a friend or family member. 4. Good news about her health. 5. Positive relationship with the medical staff.</p>
User Question Emotional Trigger	<p><i>What might have caused the woman in the image to appear irritated or angry?</i></p> <p>1. Service issue (mistake in order, long wait, problem with payment process). 2. Unpleasant environment (noise levels, cleanliness, presence of other customers). 3. Dissatisfaction with food or service. 4. Frustration or annoyance with the conversation or situation.</p>

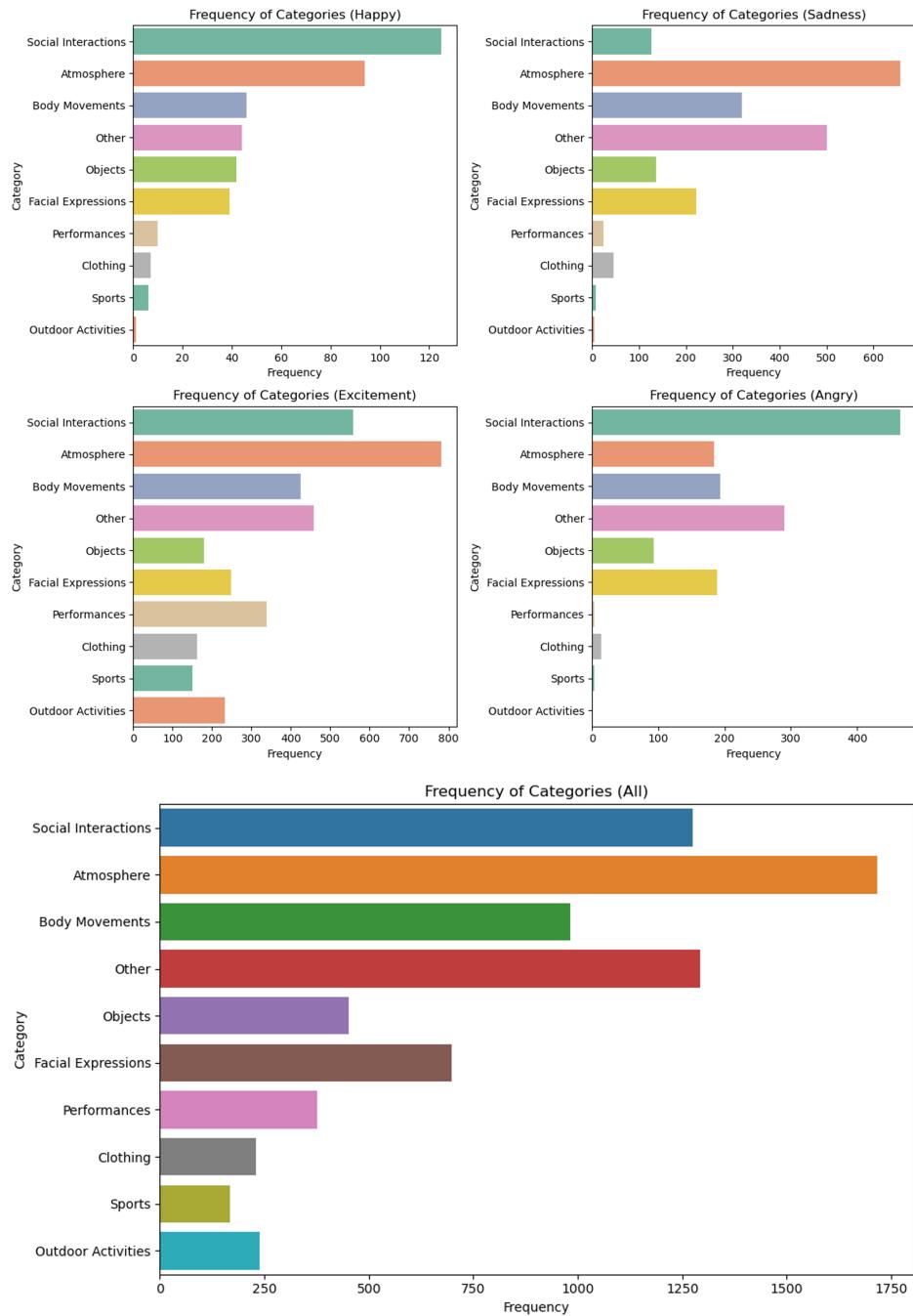


Figure 5: Visualization of the numbers of emotional triggers across different categories (Basic Emotions).

Table 9: Statistics of the Emotional Trigger Types (Basic Emotions).

Atmosphere	Social Interactions	Body Movements	Facial Expressions	Objects	Performances	Outdoor Activities	Clothing	Sports	Other
23.11%	17.17%	13.24%	9.40%	6.07%	5.06%	3.20%	3.08%	2.25%	17.41%

## A.2 Complex Emotion Comprehension Subset

Table 10: Visualization of complex emotion comprehension subset, an image is corresponded to multiple user questions.

Examples of the Complex Emotion Comprehension Subset	
User Question (1) Emotional Trigger  User Question (2) Emotional Trigger	 <p><i>Why does the kid in the background seem excited?</i>          1. Head turning back. 2. Starring at the two playing with each other on the focus. 3. Sense of motion from the event. 4. Maybe excited about the desire to join them.  <i>What do you think might have caused the kid in the background of the image to be confused?</i>          1. Head turning back. 2. Two others acting abnormally. 3. Two others each holding a stick of corn. 4. Maybe curious about the event. 5. Maybe wondering about the motivation for the abnormality.</p>
User Question (1) Emotional Trigger User Question (2) Emotional Trigger User Question (3) Emotional Trigger	 <p><i>What may caused the little girl upset?</i>          1. Crying. 2. Can not making handiwork. 3. The woman blamed her.  <i>What may caused the little girl happy?</i>          1. Crying but the women comfort her. 2. Can not making handiwork. 3. Woman help her finishing the work.  <i>What may cause the woman angry?</i>          1. The girl is not obedient. 2. The girl can't do handiwork. 3. The girl can't learn no matter how much taught. 4. Step-by-step instruction.</p>
User Question (1) Emotional Trigger User Question (2) Emotional Trigger	 <p><i>Why does the baby show the fear expression?</i>          1. The man's scary outfit. 2. Afraid of the man. 3. The man's makeup. 4. Covering mouth with hand.  <i>What make the baby surprise and happy?</i>          1. Shocking face and gesture. 2. Staring at someone. 3. Sense of unbelievable. 4. A man colored in silver on the focus. 5. Maybe shocked to see something abnormal.</p>
User Question (1) Emotional Trigger User Question (2) Emotional Trigger	 <p><i>Why does this man in the picture look exhausted and annoyed?</i>          1. Maybe lack of Sleep. 2. Closed-eyes. 3. Taking care of a young child. 4. Tired of the child. 5. Naughty child.  <i>Why does this man being enjoyment and pleasure?</i>          1. Enjoying spending time with his child. 2. Child lying in arms. 3. Satisfied with the moment. 4. Sense of company of family. 5. Engaging in playful activities.</p>

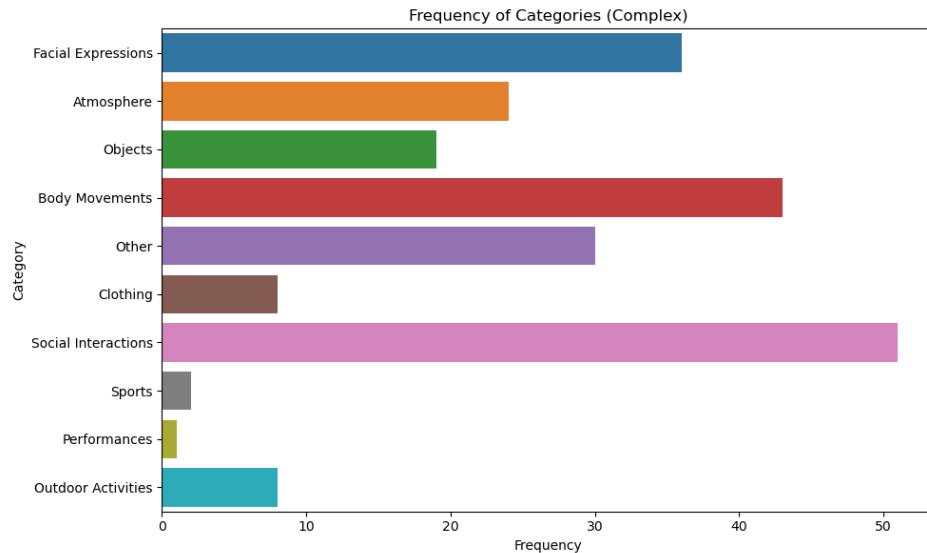


Figure 6: Visualization of the numbers of emotional triggers in the Complex Emotion Comprehension Subset.

Table 11: Statistics of the Emotional Trigger Types (Complex Emotions).

Atmosphere	Social Interactions	Body Movements	Facial Expressions	Objects	Performances	Outdoor Activities	Clothing	Sports	Other
10.81%	23.00%	19.37%	16.22%	8.55%	0.45%	3.60%	3.60%	0.9%	13.51%

## B Human in the Loop Data Cleaning Details

### B.1 Hallucinations in VLLMs

In this section, we present examples of hallucinations in Vision Large Language Models (VLLMs), along with the human-in-the-loop data cleaning process to address them. Table 12 provides examples of hallucinated emotional triggers generated by VLLMs in response to user questions. The examples illustrate instances where VLLMs generate triggers that do not accurately reflect the visual context of the images, or are not present in the image. For instance, in the first example, the VLLM hallucinates “Doing mountain biking” as the trigger for the man’s participation in the outdoor activity, despite no evidence of biking in the image. Removing these hallucinations reduces bias in our datasets introduced by VLLMs.

Table 12: Example of Hallucinations in VLLMs. Hallucinations are indicated in red, while other text is indicated in gray.

Examples of the Human Cleaning Process of Hallucinations	
User Question	<i>What might have motivated the man in the image to participate in this outdoor activity, given his gear and the environment?</i>
Emotional Trigger (Raw)	1. Determination and concentration. 2. Challenge of the race or trail. 3. Personal goals. 4. Desire to improve mountain biking skills. 5. Well-prepared gear. 6. Environmental factors (rocky slope, weather conditions). <b>7. Doing mountain biking.</b>
User Question	<i>What could have caused the man in the image to appear outraged or hostile?</i>
Emotional Trigger (Raw)	1. Holding a black bag. 2. Animated conversation or gesture. 3. Furrowed eyebrows. 4. Open mouth. 5. Wide or squinting eyes. 6. Leaning forward or gesturing with hands. <b>7. Brown couch (as a place where he typically relaxes or discusses matters)</b>
User Question	<i>What might have caused the man in the image to be angry or upset?</i>
Emotional Trigger (Raw)	1. KANO CAP ABILITY sign on the wall. 2. Feeling overwhelmed or pressured by his workload. 3. Undervalued or overworked in his professional role 4. Recent events or interactions in the workplace that have caused stress or frustration. 5. Tension and stress in his body language (posture, grip on the mug).
User Question	<i>What might have caused the man in the image to appear angry or frustrated?</i>
Emotional Trigger (Raw)	1. Disagreement with a family member. 2. Concern about a meal he is preparing. 3. Problem at work that he is thinking about while in the kitchen. 4. Serious or intense mood due to work-related issue or concern.

## B.2 Examples of Data Cleaning for Commonsense Knowledge

In this section, we provide examples of the human-in-the-loop data cleaning process to instill commonsense knowledge into the dataset. Table 13 presents instances where commonsense knowledge is incorporated to improve the emotional trigger identification. By incorporating commonsense knowledge, we enhance the accuracy and realism of emotional triggers, improving the overall quality of the dataset for emotion comprehension tasks.

Table 13: The Human in the Loop process instills Commonsense Knowledge into the dataset. Text orange represents added commonsense knowledge.

Examples of Data Cleaning for Commonsense Knowledge	
User Question	
Emotional Trigger	1. Halloween costume and bib with a pumpkin design. 2. Interaction with the person holding them up. 3. Festive atmosphere and attention from the person holding them up. 4. First Halloween experience.
User Question	
Emotional Trigger	1. A toy written "Beijing Welcome". 2. Taking a photo with Tienanmen Square. 3. First time to Beijing.
User Question	
Emotional Trigger	1. Celebratory event or milestone related to the year 2021. 2. Excitement and joy. 3. Playful or lighthearted moment shared between the man and the woman. 4. Achievement or personal milestone. 5. Festive and celebratory atmosphere.
User Question	
Emotional Trigger	1. Head turning back. 2. Starring at the two playing with each other on the focus. 3. Sense of motion from the event. 4. Maybe excited about the desire to join them.

## C Case Study of the VLLMs Emotion Comprehension Abilities

**Table 14:** Case Study of the VLLMs Emotion Comprehension Abilities, each example present the emotion comprehension process and the LLaMA-3/ChatGPT evaluation.

Examples of VLLMs Emotion Comprehension	
User Question Emotional Trigger	
Qwen-VL-Chat	<i>What factors might be contributing to the anger of the man in the image?</i> 1. Staring with wide opening eyes. 2. Fully aggressive vibe. 3. Metal claws. 4. Defending gesture. 5. Shrunk muscle.
LLaMA-3 Eval	The man in the image is holding several metal claws in his hands and making a fierce face, which suggests that he is angry. Factors that might be contributing to his anger could include a conflict or disagreement with someone, a frustrating situation, or a sense of injustice or betrayal. It's also possible that the man is simply expressing his emotions in a playful or dramatic way for fun or to entertain others... {score: 2/5}. *Correctly identified: The man's facial expression and the metal claws in his hands are correctly identified as potential triggers for his anger (points 2 and 3). *Missed: The other factors mentioned in the ground truth label, such as staring with wide opening eyes, defending gesture, and shrunk muscle, are not mentioned in the record.
ChatGPT Eval	{score: 3/5}. Explanation: The record correctly identified the metal claws and the defending gesture as potential triggers of anger. However, it missed identifying the staring with wide opening eyes, fully aggressive vibe, and shrunk muscle as additional potential triggers.
Video-LLaVA	The man in the image appears to be angry, possibly due to the presence of multiple knives in his hands. These knives could be causing him distress or discomfort, or he might be experiencing frustration or aggression. Additionally, the man's facial expression and body language suggest that he is feeling overwhelmed or threatened. It is important to note that the man's emotions may be influenced by multiple factors, such as personal circumstances, the environment, or the presence of others around him. {score: 1/5}. The record correctly identifies one potential trigger for the man's emotion, which is the presence of multiple knives in his hands. However, the other triggers mentioned, such as facial expression, body language, and personal circumstances, are not explicitly mentioned in the ground truth label. The ground truth label specifically mentions staring with wide opening eyes, metal claws, defending gesture, and shrunk muscle, which are not mentioned in the record. {score: 2/5}. Explanation: The record correctly identified the aggressive vibe and defending gesture as potential triggers for the emotion of anger. However, it missed identifying the staring with wide opening eyes, metal claws, and shrunk muscle.
MiniGPT-v2	The factors contributing to the anger of the man in the image could be his physical posture, facial expression, and hand gesture. His raised hands with his fingers interlocked, forming a claw-like shape, suggest that he is expressing aggression or hostility. His furrowed brow and intense gaze further reinforce this interpretation. Additionally, his clenched fists, which are positioned in front of his chest, indicate that he might be experiencing strong emotions such as frustration, irritation, or even rage. These elements together suggest that the man in the image may be expressing a range of negative emotions, which can contribute to his overall appearance of anger. {score: 2/5}. The record correctly identified the following potential triggers: *Raised hands with fingers interlocked, forming a claw-like shape (matches ground truth label 3: Metal claws). *Furrowed brow and intense gaze (matches the overall aggressive vibe described in ground truth label 2: Fully aggressive vibe). The record missed the following potential triggers: *Staring with wide opening eyes (not mentioned in the record). *Defending gesture (not explicitly mentioned, but the record does mention hand gesture, which could be interpreted as a defensive gesture, but it's not clear). *Shrunk muscle (not mentioned at all in the record). *Clenched fists (while mentioned, it's not explicitly linked to the emotions described in the ground truth label, which is frustration, irritation, or rage). The record provides a good analysis of the physical posture, facial expression, and hand gesture, but misses some specific details mentioned in the ground truth label. {score: 3/5}. The record correctly identified the factors of physical posture, facial expression, and hand gesture contributing to the man's anger. The raised hands forming a claw-like shape, furrowed brow, and intense gaze were correctly interpreted as signs of aggression or hostility. However, the record missed identifying the wide opening eyes and fully aggressive vibe as triggers of anger. The interpretation of defending gesture and shrunk muscle were not mentioned in the record.
Otter	Several factors might contribute to the man's angry expression in this image. He could be upset due to the intense lighting shining on his face, which might make him appear angrier than he actually is. Additionally, the man is holding a pair of large scissors, which could be potentially intimidating or making him feel uncomfortable. The close-up nature of the image might also emphasize his facial expressions, accentuating his anger. Furthermore, the context of the scene could be a photo shoot or a situation where the man might be experiencing emotions related to the situation he is in. {score: 2/5}. The record correctly identifies two potential triggers for the man's angry expression: the intense lighting shining on his face and the close-up nature of the image emphasizing his facial expressions. However, it misses three potential triggers mentioned in the ground truth label: the man's staring with wide opening eyes, the fully aggressive vibe, and the metal claws. {score: 1/5}. Explanation: The record correctly identified the potential trigger related to the man holding a pair of large scissors, which could be potentially intimidating or making him feel uncomfortable. However, it missed the other triggers mentioned in the ground truth label, such as staring with wide opening eyes, fully aggressive vibe, metal claws, defending gesture, and shrunk muscle.
LLaMA-3 Eval	
ChatGPT Eval	
ChatGPT-4V	The man in the image is showing a facial expression and body language that can be interpreted as anger or aggression. Key features contributing to this interpretation may include: 1. Eyebrows: His eyebrows are furrowed, creating wrinkles on his forehead, which is a common indicator of anger or concentration. 2. Eyes: His eyes have an intense gaze which suggests focus and potential hostility. 3. Mouth: His mouth is closed tightly, which can be a sign of tension or suppressed emotion. 4. Facial muscles: There is a general tension in his facial muscles, adding to the stern expression. 5. Stance: His body stance appears confrontational, with his arms and hands in front of him as if ready for combat. 6. Props: The metallic claws he's holding, reminiscent of the character Wolverine from the X-Men series, suggest a readiness for defense or attack, fitting with the theme of aggression or anger... {score: 4/6}. The record correctly identified the following potential triggers: *Eyebrows: Furrowed eyebrows are a common indicator of anger or concentration, which is correct. *Eyes: An intense gaze can suggest focus and potential hostility, which is correct. *Mouth: A closed mouth can be a sign of tension or suppressed emotion, which is correct. *Facial muscles: General tension in facial muscles adds to the stern expression, which is correct. *Props: The metallic claws suggest a readiness for defense or attack, fitting with the theme of aggression or anger, which is correct. The record missed the following potential triggers: *Stance: The ground truth label mentions a "defending gesture", which is not explicitly mentioned in the record. *Muscle: The ground truth label mentions "shrunk muscle", which is not mentioned in the record... {score: 3/5}. Explanation: The record correctly identified the intense gaze in the eyes, the confrontational stance, and the presence of metal claws as potential triggers of aggression or anger. However, it missed mentioning the wide opening eyes and the defending gesture as described in the ground truth label. The mention of shrunk muscles was not identified in the record.
LLaMA-3 Eval	
ChatGPT Eval	

Continue on next page

Examples of VLLMs Emotion Comprehension	
Emotional Trigger	1. Staring with wide opening eyes. 2. Fully aggressive vibe. 3. Metal claws. 4. Defending gesture. 5. Shrunk muscle.
LLaVA-1.5 (13B)	The man in the image, who is dressed in a flannel shirt and jeans and wearing a belt, appears to be angry. Several factors might be contributing to his anger. One possible factor is that he is posing with a set of fake Wolverine claws, which could indicate that he is a fan of the X-Men character Wolverine and is expressing his frustration or disappointment about the character's portrayal in movies or other media. Another factor could be that he is simply in a bad mood or experiencing personal issues that are unrelated to the claws. It is also possible that the man is posing for a humorous or satirical photo, which might be intended to provoke a reaction or make people laugh... (score: 2/5). The record correctly identifies two potential triggers for the man's anger: the fake Wolverine claws (factor 3) and the possibility that he is a fan of the X-Men character Wolverine and is expressing frustration or disappointment (factor 1). However, it misses three other potential triggers mentioned in the ground truth label: the man's staring with wide opening eyes (factor 1), the fully aggressive vibe (factor 2), and the defending gesture (factor 4).
LLaMA-3 Eval	[score: 3/5]. Explanation: The record correctly identified the metal claws and the defending gesture as potential triggers for the emotion of anger. However, it missed identifying the staring with wide opening eyes, fully aggressive vibe, and shrunk muscle as additional potential triggers.
ChatGPT Eval	[score: 3/5]. Explanation: The record correctly identified the metal claws and the defending gesture as potential triggers for the emotion of anger. However, it missed identifying the staring with wide opening eyes, fully aggressive vibe, and shrunk muscle as additional potential triggers.
LLaVA-NEXT (7B)	In the image, you see a person with a stern expression, holding a set of six stainless steel claws that resemble those of Wolverine from the "X-Men" comics and movies. The factors contributing to the anger of the man in the image could include: 1. **Stress or Discomfort**: The claws might be causing him physical discomfort or pain, which can lead to an angry expression. 2. **Frustration**: He might be frustrated with a situation or a person, and the claws could be a symbol of his anger. 3. **Emotional Exhaustion**: The act of holding the claws might be emotionally exhausting, and the man could be angry because he feels overwhelmed. 4. **Performance or Role**: The man could be an actor or a cosplayer portraying a character with anger, and the claws are part of his costume. 5. **Personal Issues**: There could be personal issues or conflicts that are causing him anger, and the claws could be a metaphor for his feelings of aggression or frustration. 6. **Artistic Expression**: The image might be staged for an artistic purpose, and the man's expression and the claws are part of a creative concept. 7. **Reaction to a Situation**: The man might be reacting to a recent event or situation that has caused him anger, and the claws are a representation of his emotional response. It's important to note that without additional context, it's difficult to pinpoint the exact cause of his anger. However, these are some plausible factors that could be contributing to his expression. (score: 1/7). My assessment is that 2 out of the 7 potential triggers identified in the record match the ground truth label. The two correct triggers are: * **Frustration**: The ground truth label mentions a fully aggressive vibe, which aligns with the idea that the man might be frustrated with a situation or a person. * **Reaction to a Situation**: The ground truth label mentions a defending gesture, which suggests that the man might be reacting to a recent event or situation that has caused him anger. The other 5 potential triggers do not directly match the ground truth label.
LLaMA-3 Eval	[score: 1/5]. Explanation: The record correctly identified the presence of metal claws, which is one potential trigger for the emotion. However, it missed the other triggers such as staring with wide opening eyes, fully aggressive vibe, defending gesture, and shrunk muscle as indicated in the ground truth label.
ChatGPT Eval	[score: 1/5]. Explanation: The record correctly identified the presence of metal claws, which is one potential trigger for the emotion. However, it missed the other triggers such as staring with wide opening eyes, fully aggressive vibe, defending gesture, and shrunk muscle as indicated in the ground truth label.
LLaVA-NEXT (13B)	The image shows a man with a stern and intense expression, holding up what appear to be claw-like objects. His facial expression, body language, and the context of the image suggest several factors that might be contributing to his anger: 1. **Personal Distress**: The man's facial expression, with furrowed brows and a tight-lipped mouth, indicates that he is experiencing significant emotional distress. This could be due to personal issues, frustration with a situation, or a perceived injustice. 2. **Physical Discomfort**: The claw-like objects he is holding could imply that he is experiencing physical pain or discomfort, which might be contributing to his anger. 3. **Perceived Threat**: The claw-like objects might be interpreted as a weapon or a defensive posture, suggesting that the man feels threatened or under attack, which could be a source of his anger. 4. **Role or Character**: If the man is portraying a character from a movie or a role-playing scenario, his anger might be part of the character's persona or the narrative of the scene. 5. **Context of the Photo**: The context of the photo, such as the setting or the people around him, could be contributing to his anger. For example, he might be in a situation where he feels cornered or unsupported. 6. **Emotional Trigger**: There could be an event or situation that has triggered his anger, which is not visible in the image... (score: 3/6). My assessment is based on the record's identification of potential triggers for the man's anger. The record correctly identifies the following triggers: 1. **Personal Distress**: The record mentions the man's facial expression, which is consistent with the ground truth label of "staring with wide opening eyes", indicating a sense of distress or alarm. 2. **Physical Discomfort**: The record mentions the claw-like objects, which is consistent with the ground truth label of "metal claws", suggesting physical discomfort or pain. 3. **Perceived Threat**: The record mentions the claw-like objects as a potential source of the man's anger, which is consistent with the ground truth label of "fully aggressive vibe" and "defending gesture", indicating a perceived threat. The record misses the following triggers: 1. **Role or Character**: The record does not explicitly mention the man's role or character, which is mentioned in the ground truth label. 2. **Context of the Photo**: The record does not provide specific context about the photo, which is mentioned in the ground truth label. 3. **Emotional Trigger**: The record does not provide a specific event or situation that triggered the man's anger, which is mentioned in the ground truth label. 4. **Shrunk muscle**: The record does not mention the man's physical posture or muscle tension, which is mentioned in the ground truth label...
LLaMA-3 Eval	[score: 2/5]. Explanation: The record correctly identified the potential triggers of "Physical Discomfort" and "Perceived Threat" based on the man's expression and the objects he is holding. However, it missed identifying triggers related to "Staring with wide opening eyes," "Fully aggressive vibe," and "Defending gesture" as described in the ground truth label.
ChatGPT Eval	[score: 2/5]. Explanation: The record correctly identified the potential triggers of "Physical Discomfort" and "Perceived Threat" based on the man's expression and the objects he is holding. However, it missed identifying triggers related to "Staring with wide opening eyes," "Fully aggressive vibe," and "Defending gesture" as described in the ground truth label.
ChatGPT-4o	Based on the image, several factors might be contributing to the man's expression of anger or intensity: 1. **Body Language**: The man is clutching blades that are extending from his hands, which might symbolize aggression or readiness for battle, contributing to a perception of anger or intensity. 2. **Facial Expression**: His furrowed brows, intense eye contact, and tight lips create an expression commonly associated with anger or determination. 3. **Costume and Pose**: The attire and the claws suggest he might be cosplaying a character known for aggression or conflict, which might be contributing to his intense demeanor. These elements together create an impression of anger or a serious mood, likely for dramatic effect, possibly related to the character he is portraying.
LLaMA-3 Eval	[score: 3/5]. My assessment is based on the record's identification of potential triggers contributing to the man's expression of anger or intensity. Here's a breakdown of the correct and missed triggers: Correctly identified triggers: 1. **Body Language**: The record correctly notes the man's clutching of blades as a symbol of aggression or readiness for battle, which contributes to a perception of anger or intensity. 2. **Facial Expression**: The record accurately identifies the man's furrowed brows, intense eye contact, and tight lips as an expression commonly associated with anger or determination. 3. **Costume and Pose**: The record correctly notes the attire and claws suggesting the man might be cosplaying a character known for aggression or conflict, which contributes to his intense demeanor. Missed triggers: 1. **Staring with wide opening eyes**: The record does not mention the man's eye expression, which is a key aspect of his intense demeanor. 2. **Fully aggressive vibe**: The record does not explicitly state the man's overall aggressive vibe, although it does mention the costume and pose suggesting aggression. 3. **Defending gesture**: The record does not mention the man's defending gesture, which is an important aspect of his intense demeanor. 4. **Shrunk muscle**: The record does not mention the man's muscle tone, which is not explicitly stated in the ground truth label. Overall, the record correctly identifies three out of five potential triggers contributing to the man's expression of anger or intensity.
ChatGPT Eval	[score: 3/5]. Explanation: The record correctly identified the factors related to body language, facial expression, and costume that contribute to the perception of anger or intensity. However, it missed identifying the specific details of staring with wide opening eyes and defending gesture as mentioned in the ground truth label.

Continue on next page

#### Examples of VLLMs Emotion Comprehension

Emotional Trigger	1. Staring with wide opening eyes. 2. Fully aggressive vibe. 3. Metal claws. 4. Defending gesture. 5. Shrunk muscle.
LLaVA-NEXT (34B)	The image you've provided shows a person with claws protruding from their hands, which is a characteristic feature of the fictional character Wolverine from the X-Men series. The person is also making a stern facial expression, which could be interpreted as a sign of anger or aggression. Factors that might contribute to the anger or intense expression in the image could include: 1. Role Playing: The person might be a fan of the character Wolverine and could be role-playing or imitating the character's typical intense demeanor. 2. Costume or Prop: The claws could be part of a costume or prop for a cosplay event or a themed party, where the person is trying to embody the character's persona. 3. Performance: The person could be an actor or performer engaging in a theatrical or filmic portrayal of the character, and the expression is part of the character's emotional state in that scene. 4. Personal Expression: The person might be expressing their own emotions, perhaps frustration, anger, or determination, and is using the claws and intense expression as a metaphor or a way to convey their feelings. 5. Artistic Choice: The image could be a piece of art or a still from a video where the artist or performer has chosen to portray the character in a specific emotional state for creative or storytelling purposes. It's important to note that without additional context, it's difficult to determine the exact reason for the person's expression. The factors listed above are speculative and based on the visual cues present in the image.
LLaMA-3 Eval	{score: 2/5}. My assessment is based on the record provided, which attempts to understand the emotion in the image by identifying potential triggers. The record correctly identifies two factors that might contribute to the anger or intense expression in the image: 1. Role Playing: The record correctly notes that the person might be a fan of the character Wolverine and could be role-playing or imitating the character's typical intense demeanor, which aligns with the ground truth label's description of a "fully aggressive vibe". 2. Performance: The record also correctly suggests that the person could be an actor or performer engaging in a theatrical or filmic portrayal of the character, which is consistent with the ground truth label's description of a "defending gesture". However, the record misses three potential triggers: 1. Staring with wide opening eyes: The record does not mention this specific visual cue, which is present in the ground truth label. 2. Metal claws: While the record correctly notes the presence of metal claws, it does not explicitly link them to the emotion or expression. 3. Shrunk muscle: The record does not mention this specific detail, which is present in the ground truth label. Overall, the record provides some insightful observations, but could benefit from more specific and detailed analysis to accurately capture the emotion and expression in the image.
ChatGPT Eval	{score: 1/5}. Explanation: The record correctly identified the presence of metal claws in the image, which aligns with the ground truth label. However, it missed identifying the other emotional triggers such as staring with wide opening eyes, fully aggressive vibe, defending gesture, and shrunk muscle. The analysis focused more on speculative factors related to role-playing, costume, performance, personal expression, and artistic choice rather than accurately pinpointing the specific emotional cues present in the image.
Qwen-vl-plus	Based on the image alone, it is difficult to determine with certainty why the man appears angry or aggressive. However, there could be several reasons for his expression: - He may have had a frustrating experience that triggered negative emotions. - The pose and stance he has adopted suggest confidence and assertiveness which can also come across as aggression. It's important to note that interpreting human emotions from images can often lead to misinterpretations due to cultural differences, personal biases, and lack of context. Therefore, any conclusions drawn should not be taken too seriously without additional information about the situation.
LLaMA-3 Eval	{score: 0/5}. The record does not identify any potential triggers correctly. The record provides a general analysis of the image, mentioning possible reasons for the man's expression, but does not specifically identify any triggers that match the ground truth label. The ground truth label includes specific details such as staring with wide opening eyes, metal claws, and shrunk muscle, which are not mentioned in the record.
ChatGPT Eval	{score: 0/5}. Explanation: The record of understanding did not correctly identify any of the potential triggers mentioned in the ground truth label. The reasons provided for the man's expression were more general and did not align with the specific details mentioned in the ground truth.
LLaMA-3 (8B)	Based on the image alone, it is difficult to determine with certainty why the man appears angry or aggressive. However, there could be several reasons for his expression: - He may have had a frustrating experience that triggered negative emotions. - The pose and stance he has adopted suggest confidence and assertiveness which can also come across as aggression. It's important to note that interpreting human emotions from images can often lead to misinterpretations due to cultural differences, personal biases, and lack of context. Therefore, any conclusions drawn should not be taken too seriously without additional information about the situation.
LLaMA-3 Eval	{score: 0/5}. The record does not identify any potential triggers correctly. The record provides a general analysis of the image, mentioning possible reasons for the man's expression, but does not specifically identify any triggers that match the ground truth label. The ground truth label includes specific details such as staring with wide opening eyes, metal claws, and shrunk muscle, which are not mentioned in the record.
ChatGPT Eval	{score: 0/5}. Explanation: The record of understanding did not correctly identify any of the potential triggers mentioned in the ground truth label. The reasons provided for the man's expression were more general and did not align with the specific details mentioned in the ground truth.