Facial Affective Behavior Analysis with Instruction Tuning

Yifan Li¹, Anh Dao¹, Wentao Bao¹, Zhen Tan², Tianlong Chen^{3,4,5}, Huan Liu², and Yu Kong¹

Michigan State University, {liyifa11, anhdao, baowenta, yukong}@msu.edu
Arizona State University, {ztan36, huanliu}@asu.edu
³ University of North Carolina at Chapel Hill
⁴ Massachusetts Institute of Technology, tianlong@mit.edu
⁵ Harvard University

Abstract. Facial affective behavior analysis (FABA) is crucial for understanding human mental states from images. However, traditional approaches primarily deploy models to discriminate among discrete emotion categories, and lack the fine granularity and reasoning capability for complex facial behaviors. The advent of Multi-modal Large Language Models (MLLMs) has been proven successful in general visual understanding tasks. However, directly harnessing MLLMs for FABA is challenging due to the scarcity of datasets and benchmarks, neglecting facial prior knowledge, and low training efficiency. To address these challenges, we introduce (i) an instruction-following dataset for two FABA tasks, i.e., facial emotion and action unit recognition, (ii) a benchmark FABA-Bench with a new metric considering both recognition and generation ability, and (iii) a new MLLM EmoLA as a strong baseline to the community. Our initiative on the dataset and benchmarks reveal the nature and rationale of facial affective behaviors, i.e., finegrained facial movement, interpretability, and reasoning. Moreover, to build an effective and efficient FABA MLLM, we introduce a facial prior expert module with face structure knowledge and a low-rank adaptation module into pre-trained MLLM. We conduct extensive experiments on FABA-Bench and four commonly-used FABA datasets. The results demonstrate that the proposed facial prior expert can boost the performance and EmoLA achieves the best results on our FABA-Bench. On commonly-used FABA datasets, EmoLA is competitive rivaling taskspecific state-of-the-art models. The dataset and codes are available: https://johnx69.github.io/FABA/.

Keywords: Facial affective behavior analysis \cdot Emotion recognition \cdot AU recognition \cdot Multi-modal large language model

1 Introduction

Facial affective behavior analysis (FABA) [46,83] such as facial emotion recognition (FER) [55,75], and action unit recognition (AUR) [26], aims to recognize

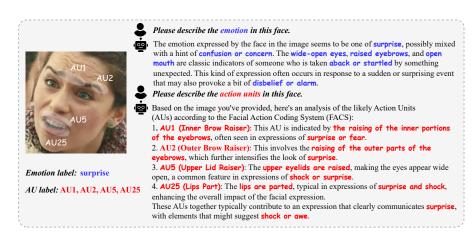


Fig. 1: An illustration of FABA-Instruct annotations. FABA-Instruct can provide fine-grained emotion and AU descriptions, which not only include the reasoning process about the facial movements but also present the inference to the emotion. Furthermore, compared to traditional category labels, FABA-Instruct has more abundant expressions to describe complex, nuanced, exaggerated, and undefined affective behaviors.

facial expressions and movements, which are critical to understanding an individual's emotional states and intentions [32]. FABA has emerged as a burgeoning field with potential across multiple domains. For example, in psychology, FABA can aid therapists by offering real-time insights into a patient's unspoken emotions through facial expression analysis, thereby enhancing therapeutic outcomes [93]. In education, it improves e-learning experiences by adjusting content delivery based on students' facial cues, indicating engagement or confusion [98].

Despite promising progress being made, most of the existing FABA approaches [64, 87, 101, 116] are based on discriminative models, which treat FER or AUR as a multi-class or multi-label classification task. Such approaches tend to induce shortcomings like coarse-grained emotional descriptions, inability to describe complex emotions, and lack of reasoning ability. Those limitations hinder the applications of FABA, for instance, in providing nuanced feedback to therapists about a patient's emotional nuances in psychology, or in accurately adapting educational content based on subtle student reactions in e-learning environments. To counteract these drawbacks, we are motivated by the success of recent multi-modal large language models (MLLMs) [120, 136], because of their evidenced ability to describe and reason over fine-grained and complex visual cues by instruction tuning after large-scale pre-training [71]. In practice, MLLMs transform the discriminative task into a sequence-to-sequence generative one [13, 16] based on large language models (LLMs) [43, 114, 115]. MLLMs have shown great capability on various visual understanding tasks such as visual question answering [70,71,152], captioning [52,53], grounding [11,133], segmentation [48], etc.

However, there exist three major challenges in FABA tasks when deploying MLLMs. (1) There is no suitable FABA dataset for MLLMs to perform instruction tuning. Existing FABA datasets have either coarse-grained annotations [28, 47, 54, 56, 84, 142, 143] or limited emotion descriptions [74] for video clips, since it's labor-intensive and expensive to manually annotate large-scale fine-grained FABA descriptions especially for AUR task. (2) There is an increasing number of MLLMs, how to select suitable MLLMs for FABA remains unknown. Existing metrics for evaluating MLLMs are language-oriented, without specific consideration of language usage in FABA tasks. (3) Furthermore, image features from vision encoders like CLIP [94] of current MLLMs struggle to capture facial structure information such as facial landmarks, leaving the impact of facial priors on FABA tasks unexplored. Fine-tuning the entire model hosting billions of parameters for these features leads to prohibitive computational costs.

To solve these challenges, we propose an instruction-following FABA dataset "FABA-Instruct". It includes 19K in-the-wild aligned face images with 30K fine-grained emotion and AU annotations using GPT4V enabling instruction tuning (Fig. 1). Based on this dataset, we propose a new benchmark "FABA-Bench" for evaluating both the visual recognition and text generation performance of various MLLMs on FABA tasks. Moreover, we introduce an efficient MLLM "EmoLA" for FABA tasks by incorporating a low-rank adaptation method and a facial prior expert to a pre-trained MLLM like LLaVA-1.5 [70]. Specifically, to obtain the facial prior knowledge, we utilize a pre-trained face alignment encoder to extract the facial landmark features, which are complementary to the vision encoder. To mitigate the computational cost, the LoRA [39] method is adopted in training such that only the parameter residual is learned through low-rank matrices. In this paper, we take the earliest trial of the FABA tasks with instruction tuning over MLLMs, which shed light on both FABA and MLLMs research community. Our contributions are summarized in three-folds:

- Instruction-following FABA data. To our best knowledge, this is the
 first FABA dataset that enables instruction tuning. It reveals new aspects of
 FABA research topics and it will continuously bring the benefits of MLLMs
 to the FABA community.
- *Instruction-following FABA benchmark*. To evaluate the recognition and reasoning ability of different models on instruction-following FABA tasks, we introduce the FABA-Bench benchmark with a unique metric, *i.e.*, REGE, allowing for both recognition and generation capability.
- MLLM-based FABA architecture. To efficiently train on FABA tasks and utilize the facial prior knowledge, we introduce the EmoLA model which involves tuning LoRA on a pre-trained MLLM and incorporating a facial prior expert. We demonstrate the effectiveness of EmoLA on FABA-Instruct and four traditional FABA datasets. The results show that EmoLA achieves the best performance on FABA-Instruct and SOTA-comparable or even better results on the traditional FABA datasets.

4

2 Related work

2.1 Facial affective behavior analysis

Psychology perspective. According to psychology research [32], two mechanisms can be utilized to model facial affective behaviors, i.e., emotion categories, and dimensional theory. According to Izard [41] and Ekman [27], the basic emotions can be categorized into one of the several prototypes. Ekman [26] also proposed to decompose the macroscopic affective behaviors into fine-grained Action Units (AUs) from the anatomical perspective. However, such discrete emotion representations may be not sufficient to capture complex and fine-grained emotions. Dimensional theory describes the continuous emotions from the Euclidean space perspective. Russel adopted two dimensions: arousal and valence [97] to represent pleasantness and degree of feelings. Although the dimensional theory allows for a more nuanced understanding of emotions, it is challenging to measure and recognize for humans. We argue that human-generated descriptions offer a superior way of characterizing facial affective behaviors (see Fig. 3 and Fig. 4), which not only capture the complexity and subtlety of affective behaviors but are also more accessible and quantifiable for humans.

Methodology perspective. To better recognize the affective behaviors, current research mainly focuses on deep-learning-based techniques. These approaches can be categorized into three streams according to the task types, *i.e.*, facial emotion recognition [8,63,106,121,139,146], action unit recognition [6,19,42,61, 104,113,130,132,141] and valance-arousal regression [62,85,100,149]. Existing methods focus on capturing the fine-grained facial movement via attention mechanism [116,121,122,125,128,145], improving generalization ability by introducing auxiliary information like facial landmarks [87,113] or extra data [61,86], exploring the relationship among emotions or AUs [51,102,147], exploiting pre-trained model like self-supervised learning [6,62,65,108,124], and probing semantic information of affective behaviors [130] using CLIP [25,94]. However, these methods are mainly discriminative-based, which fail to generate fine-grained descriptions. By contrast, our method can generate detailed descriptions based on the prior knowledge from MLLM and the facial prior expert.

2.2 Multi-modal LLMs and efficient LLM adaptation

Multi-modal LLMs. Multi-modal LLMs are getting popular in multi-modal content understanding [11, 24, 67, 117, 123]. They are built on top of LLMs [2, 15, 17, 43, 114, 140], and transform visual (videos and images) and text data into a sequence of tokens as input, resulting in generative modeling of downstream multi-modal understanding tasks by next token prediction. Specific to the image-based MLLMs, image tokens are typically encoded by CLIP vision transformer [25, 94]. Then, one of the major challenges in MLLM is how to project image token features into the language domain to better utilize the instruction-following capability of LLMs. In literature, Flamingo [1] is an early work that uses cross-attention to build the interaction between images and text tokens.

Table 1: Existing FABA datasets.

Datasets	\mathbf{AU}	Emotion	Annotation	In-the wild
RAF-DB [56]	×	1	category	×
CK+ [78]	×	/	category	×
MMI [90]	×	/	category	×
SFEW [44]	×	/	category	✓
AffectNet [85]	×	/	category	/
MAFW [74]	×	/	category & short text	✓
DFEW [44]	×	/	category	✓
FERV39K [118]	×	/	category	✓
FER2013 [31]	×	/	category	✓
DISFA [84]	1	×	category	×
GFT [30]	1	×	category	×
CASME-II [129]	/	/	category	×
BP4D [143]	1	/	category	x
EmotioNet [28]	1	✓	category	✓

FABA-Instruct

 $\begin{array}{c} {\rm category} \\ {\rm category} \\ {\rm instruction} \ \& \ {\rm description} \end{array}$

Table 2: FABA-Instruct statistics.

Statistics	Value
Total images	19877
Emotion training samples	19474
Emotion testing samples	403
Emotion description average length	50.47
AU training samples	15838
AU testing samples	325
AU description average length	207.35
Face images GPT-4V Training Test FABA-Instruct	Annotator

Fig. 2: FABA-Instruct annotation.

This design is followed by recent MLLMs [3, 14, 50, 107]. Based on the cross-attention mechanism, Q-Former [52] was proved to be a superior visual-text projector and inspires recent line of MLLM research [7, 20, 59, 152]. Recently, in contrast to the attention-based projector, LLaVA series [70,71] propose to use a simple MLP as the projector and resort to the LLMs to handle the visual-text interaction. In this paper, based on LLaVA-1.5 [70], we developed an MLLM for affective behavior analysis and benchmarked multiple MLLM baselines.

Efficient LLM adaptation. Though most existing multimodal LLMs show impressive performance on general visual-language tasks, how to adapt them to downstream applications is still challenging [126], especially with limited instructional annotation data and computing resources. This raises the surge of research on parameter-efficient fine-tuning (PEFT) [38], which avoids fine-tuning the entire LLM in training. According to [126], existing PEFT methods can be categorized into addictive [35,38,49,58,68,69,80,92,96], partial [33,109,127,135], reparametrized [22,39,73,137], hybrid [9,35,45,82], and unified fine-tuning [37, 103,119,134]. In this paper, we are interested in the re-parameterized fine-tuning as it gains increasing research attention. In this line, LoRA [39] is a pioneering work, which is inspired by the fact that modal weight adaptation has a low intrinsic rank. In training, LoRA keeps most LLM parameters frozen while only optimizing the low-rank factorized matrices of dense layers' residual. In our work, we empirically found the effectiveness of LoRA on the adaptation of MLLMs toward affective behavior analysis tasks.

3 Dataset and Benchmark

3.1 Instruction-following FABA dataset

In this section, we will present our instruction-following FABA dataset **FABA-Instruct** with two different FABA tasks, *i.e.*, emotion recognition and AU recognition. Our dataset is different from existing FABA datasets w.r.t. tasks, annotation types, and the variety of images, which is presented in Tab. 1. This

Combination of emotions The emotion depicted on the face in the image appears to be a complex one, often characterized as a mixture of sodness and discomfort, potentially with an element of a forced smile. The downturned eyebrows, elevated checks, and the slight pulling up of the lip corners suggest a pained expression that could be associated with crying or distress, yet there is an attempt to smile. This could indicate a social smile or an attempt to mask one's true feelings. Degree of emotion The person in the image appears to have a neutral expression with a slight hunt of a smile. There is released composure around the eyes and mouth, suggesting a calm or content state. The eyes are engaged and there is a mild pleasantness to the overall demeanor.

Fig. 3: Emotion description analysis. Emotion descriptions can express comprehensive emotion types like compositional emotions, exaggerated emotions, the degree of emotions, and undefined emotions, etc. In contrast, emotion categories struggle to capture such complex and nuanced emotional states.

table illustrates that most of previous FABA datasets are either task-specific [30,31,44,56,74,78,84,85,90,118,143], or laboratory collected [30,78,84,90,143]. More importantly, none of these FABA datasets have instruction-following annotations. FABA-Instruct is the inaugural dataset to offer instruction-following annotations of both AU and emotion, specifically for in-the-wild face images.

Data construction. The overall annotation pipeline is shown in Fig. 2. Specifically, we use 100 carefully designed templates as the instructions for querying GPT-4V [76, 110] on the emotions and AUs in the face. For instance, the templates for querying emotion are like "What is the emotion in this face?", "What are the action units present in this face?". More details about these templates are in the Appendix. FABA-Instruct statistics are shown in Tab. 2. For the face images, we randomly sample 19,474 and 403 face images from the training and testing set of AffectNet [85] as the training and testing face images, respectively. AffectNet is a large-scale in-the-wild facial expression database, which crawls from the Internet by querying three search engines. We align and crop these face images using the Dlib library.

However, some of the annotations for these face images cannot be obtained due to the low resolution and occlusion issues, especially for the AU annotations. As a result, after filtering out the useless annotations, we obtain 19,474 and 15,835 instructions in terms of emotions and AUs, respectively. Moreover, since some of the annotations are inaccurate or inexact, we carefully check and revise each annotation for the test sets of two tasks.

Emotion description analysis. Existing emotion datasets like RAF-DB [56], AffectNet [85], CK+ [78], FER2013 [31], MMI [90], SFEW [23], mainly adopt the seven emotion categories, *i.e.*, happiness, sadness, anger, fear, disgust, surprise, neutral. However, we argue that classifying emotions into one of the several discrete emotion categories is limited in practice since emotions are

⁶ http://dlib.net/

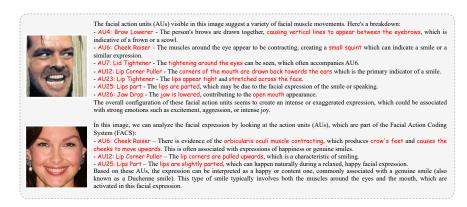


Fig. 4: AU description analysis. AU descriptions give not only the AU labels, but also provide explanations on the cause (which muscle movement) and effect (which emotion it will lead to) w.r.t. each AU, and the relationship between current AU and other AUs or emotions.

subjective [4], complex [26], continuous [72], and contextual [34]. For instance, as shown in Fig. 3, it is inaccurate to classify the compositional expression, sadness with a forced smile, into one of the emotion categories. For some exaggerated emotions, e.g., making a mouth, the actual emotion is unknown from the face image and it is not reasonable to classify the emotion according to its superficial facial movements. Also, discrete emotions cannot express the degree of emotion, and it is inexact to classify the girl's face as either happy or neutral. Furthermore, since emotions are complex, basic emotions cannot cover undefined emotions like worry, skeptical, etc. By contrast, the emotion descriptions can address these issues due to their expressive nature.

Action Unit description analysis. The annotations from traditional AUR datasets, e.g., BP4D [143], DISFA [84], EmotioNet [28], GFT [30], mainly adopt a string of binary vector to denote whether each AU is activated or not. However, the representation capability of this way is limited. It cannot indicate the degree of AU's activation and cannot provide any explanations and analysis for the prediction. For instance, in Fig. 4, merely stating that AU6 (cheek raiser) is activated cannot depict the degree of cheek muscle raising. In contrast, the AU descriptions can provide more reasoning cues like "small squint" or "crow's feet", which can capture the activated degree of AU6. Furthermore, descriptions can also show the inference ability by providing the relationship between current AU and other AUs or emotions, e.g., "indicative of a frown or a scowl", "associated with strong emotions". Such descriptions can not only provide more nuances but also improve the interpretability [12,60] of the model.

3.2 Instruction-following FABA benchmark

Evaluation. Given the fact that our FABA-Instruct uses free-form textual descriptions to represent emotions and AUs, its model evaluation stands distinct



Fig. 5: The synonyms of emotions for classifying the text.

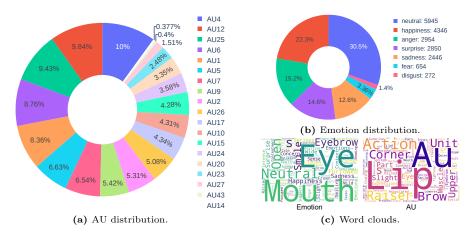


Fig. 6: The distribution of FABA-Instruct on AU (a) and emotion (b) tasks, and the word clouds (c). We extract the emotion labels using the synonyms of each emotion.

from the traditional FABA tasks and Natural Language Generation (NLG) [29] tasks. Specifically, for one thing, the traditional NLG metrics such as BLEU [91] or ROUGE [66] scores mainly concentrate on the coherence and fluency of the generative text, ignoring the FABA relevant consideration in the evaluation. For another, existing FABA metrics, e.g., accuracy or F1 score, focus only on the recognition performance of the model, lacking the evaluation on textual aspects of the model such as the reasoning and explanation. To compensate for the drawbacks of these two metrics, we introduce a new metric REGE to evaluate the **RE**cognition and **GE**neration performance of models on our FABA-Instruct.

Our REGE score is defined to consider both text generation and image recognition aspects of FABA models. For text generation, we choose to use the ROUGE score, which is a generally used metric in NLG by comparing the overlap of n-grams, word sequences, and word pairs between the generated texts and the reference ones. For recognition, we adopt the recognition accuracy for the multi-class performance of facial emotion recognition, and the F1 score for the multi-label performance of AU recognition. Denote the recognition performance as S_{re} and the generation metric as S_{ge} , our REGE metric is computed by taking their sum: $S_{rege} = S_{re} + S_{ge}$.

Calculation of S_{re} . For the FER task, we propose to classify a face image into one of the aforementioned seven basic emotion categories. In order to classify the

description, we manually select some synonyms from the training descriptions for each emotion category (see Fig. 5), and the emotion distribution of the training descriptions is shown in Fig. 6b. In practice, since the descriptions incorporate negative sentences, we first deleted these sentences. After that, we count the frequency of emotion synonyms in the sentences and treat the emotion with the highest count as the emotion label for this description. After obtaining these emotion labels of texts, we can calculate the accuracy score to evaluate the recognition performance on the FER task.

Similarly, for the AUR task, we propose to classify a face image into one or multiple AU categories. We show the distribution of AUs in Fig. 6a. Following existing literature [28, 143], we choose 12 AUs for evaluation. Subsequently, we can calculate the F1 scores to represent the recognition performance for AUR.

4 EmoLA: An Instruction-tuned MLLM for FABA

In this Section, we introduce a novel MLLM designed for FABA tasks. The overall framework of EmoLA is illustrated in Fig. 7. Its important components include two image experts (a visual and a facial prior expert), a language expert (a tokenizer with the word embedding) and a language decoder (LLM) with a LoRA module. Specifically, a face image X_V is encoded by a visual expert known as a pre-trained CLIP-L/14 [94] with a two-layer Multi-Layer Perceptron (MLP), which generates the visual embedding tokens H_v . Similarly, for the input instruction X_Q , the language expert can provide the language tokens H_q .

Notably, the visual tokens H_v may fail to capture the facial structure information since CLIP is trained with general image-text pairs rather than FABA datasets, which may make the visual expert focus more on general semantic information and overlook task-specific

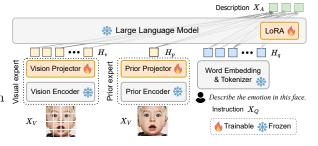


Fig. 7: EmoLA architecture.

features. Hence, we suggest employing an extra facial prior encoder $f_p(\cdot)$, trained on facial-related datasets, to better capture the facial prior knowledge and enhance the recognition ability for face images. Specifically, we adopt a pretrained facial landmark detector from Insightface⁷ to extract the landmark feature. Other face priors, such as recognition features [21], parsing, etc., can also be considered. We leave exploring these additional priors as future work for EmoLA. Therefore, the prior feature extracted by f_p can be expressed as: $Z_P = f_p(X_V)$. After obtaining Z_P , we utilize the MLP $g_\theta(\cdot)$ to project the facial prior feature Z_P to the token embedding space:

$$H_p = g_\theta(Z_P),\tag{1}$$

⁷ https://github.com/deepinsight/insightface/tree/master

where H_p is the facial prior token. This token can provide prior knowledge of face structure, which may be ignored by the visual expert.

After obtaining the visual embedding tokens H_v , facial prior token H_p and the language embedding tokens H_q , we concatenate them together and feed them into the LLM decoder. Here we use Vicuna [15] as the LLM decoder. As shown in Fig. 7, the visual encoder, prior encoder, word embedding and the LLM decoder are frozen. Except for the prior encoder $f_p(\cdot)$, the initial weights of the other modules come from a pretrained MLLM. In this paper, we utilize the LLaVA-1.5 [70] as the backbone model. In order to train efficiently without tuning the entire MLLM, we propose to tune an extra LoRA module $h_{\phi}(\cdot)$, a visual projector $h_{\gamma}(\cdot)$, and the prior projector $g_{\theta}(\cdot)$. Compared to finetuning the entire LLM, tuning LoRA reduces both memory and computation costs during training without inducing additional expenses for inference. Our experiments demonstrate the effectiveness of such a design in Tab. 7. As a result, the overall parameters to be optimized are $\Theta = \{\theta, \gamma, \phi\}$. We optimize these parameters following the auto-regressive way, and the likelihood of the generated description X_A conditioned on image X_V , prior facial feature Z_P and instructions X_Q is given by:

$$p(X_A|X_V, Z_P, X_Q) = \prod_{i=1}^{L} p_{\Theta}(x_i|X_V, Z_P, X_Q, X_{A, < i}),$$
(2)

where L is the length of the token sequence, x_i is the current token that needs to be predicted. $X_{A,< i}$ are the previous answer tokens.

5 Experiments

Implementation details. We initialize all the frozen weights of EmoLA with LLaVA-1.5 7b [70], and we only tune the prior projector and LoRA during the training stage. We train our EmoLA for one epoch optimized by AdamW with an initial learning rate of 1e-4 for all the datasets. The rank of LoRA is set to 128. We conduct all the experiments on 8 A6000 GPUs.

Database and protocols. We perform experiments on four traditional FER and AUR datasets, including one emotion

Table 3: Comparison on RAF-DB.

Methods	Accuracy (%)
RAN [116]	86.9
MA-Net [150]	87.22
EfficientFace [151]	88.36
RUL [144]	88.98
DAN [122]	89.70
EAC [145]	90.35
APViT [128]	91.98
POSTER [81]	92.05
EmoLA (Ours.)	92.05

dataset (RAF-DB [56]) and three AU datasets (BP4D [143], DISFA [84] and GFT [30]). We turn the annotations of these datasets into instruction-following ones by adding instructions. Instead of predicting the class index or binary vectors in discriminative models, our EmoLA directly outputs the corresponding emotions (i.e., "Happy") or AU labels (i.e., "AU1, AU4"). We adopt accuracy and F1 score to evaluate the recognition performance for FER and AUR tasks, respectively. More details about these datasets can be found in the Appendix. For BP4D and DISFA, following [61,87,101], we also perform a subject-exclusive

Method/AU 2 25 4 6 9 12 26 Avg. DRML [148] 17.3 17.7 37.4 29.0 10.7 37.7 38.5 20.1 26.7 EAC-Net [57] 41.5 26.4 66.450.7 80.5 89.3 88.9 15.6 48.5 DSIN [18] 42.439.0 68.1 28.6 46.8 70.8 90.4 42.2 53.6 SRERL [51] 45.7 47.8 59.647.1 73.5 84.3 43.6 55.9 45.6LP-Net [87] 29 9 24.772.746.8 496 72.993.8 65.056.9 CMS [99] 40.244.3 53.257.1 50.3 73.581.1 59.7 57.4 ARL [102] 43.942.163.6 41.8 40.076.295.266.8 58.7 SEV-Net [130] 55.3 53.161.553.6 38.271.6 95.7 41.5 58.8 HMP-PS [105] 38.065.250.950.876.093.345.967.661.0ATCM [42] 46.148.6 72.8 56.750.0 72.1 90.8 55.4 61.5 ReCoT [61] 51.336.266.8 50.152.478.8 95.369.762.6JÂA-Net [101] 62.4 60.7 67.1 41.1 45.1 73.5 90.9 67.4 63.5 PIAP [113] 63.8 50.2 51.8 71.9 50.6 54.5 79.7 94.1 57.2 GraphAU (R-50) [79] 54.6 47.172.9 54.0 55.7 76.7 91.1 53.063.1EmoLA (Ours.) 56.9 55.2 43.1 91.6

Table 4: F1 score (in %) of 8 AUs on DISFA. All the results are from original papers.

3-fold cross-validation. For GFT and RAF-FB, we train and test according to the original dataset division protocol.

We also perform experiments on our FABA-Instruct dataset w.r.t. FER and AUR tasks. We compare with other MLLMs using our SEGE metric (Sec. 3.2). We train and test all the models according to our dataset division protocol.

5.1 Comparison on traditional FER and AUR datasets

Comparison on traditional FER datasets. We conduct a comparative experiment with the latest state-of-the-art (SOTA) methods on RAF-DB, as presented in Tab. 3. Our EmoLA achieves the best results compared with previous SOTA methods. The results demonstrate the significant potential of MLLMs in addressing the FER problem. It is worth noting that most of these methods are specifically tailored for FER tasks, while our EmoLA is easy to adapt to other tasks (e.g., AUR) due to the high flexibility brought by instruction tuning.

Comparison on traditional AUR datasets. Tab. 5, Tab. 4 and Tab. 6 show the comparison results of EmoLA with other SOTA methods on AUR datasets. It can be observed that EmoLA outperforms all the other SOTA methods on DISFA and GFT datasets, achieving a 1.3% improvement over PIAP [113] on DISFA, and a 3.5% increase over EmoCo [108] on GFT. Additionally, it closely competes with ReCoT on BP4D, with a marginal gap of only 0.6%. The gap stems from the benefits of the consistency regularization and co-training in ReCoT to the BP4D dataset. We can observe that EmoLA excels in multi-label classification tasks, potentially due to its strategy of exclusively predicting the positive labels, thereby mitigating the imbalanced issue caused by negative labels.

5.2 Comparison on FABA-Bench

We compare with current MLLMs, *i.e.*, LLaVA-1.5 [70], MiniGPT4-V2 [10], Shikra [11], and mPLUG-Owl2 [131], on our FABA-Bench in Tab. 7. These

Table 5: F1 score (in %) of 12 AUs on BP4D. The results with * are taken from [61]. All the other results are taken directly from their original papers.

Method/AU	1	2	4	6	7	10	12	14	15	17	23	24	Avg.
DRML [148]	36.4	41.8	43.0	55.0	67.0	66.3	65.8	54.1	33.2	48.0	31.7	30.0	48.3
EAC-Net [57]	39.0	35.2	48.6	76.1	72.9	81.9	86.2	58.8	37.5	59.1	35.9	35.8	55.9
DSIN [18]	51.7	40.4	56.0	76.1	73.5	79.9	85.4	62.7	37.3	62.9	38.8	41.6	58.9
CMS [99]	49.1	44.1	50.3	79.2	74.7	80.9	88.3	63.9	44.4	60.3	41.4	51.2	60.6
LP-Net [51]	43.4	38.0	54.2	77.1	76.7	83.8	87.2	63.3	45.3	60.5	48.1	54.2	61.0
ARL [102]	45.8	39.8	55.1	75.7	77.2	82.3	86.6	58.8	47.6	62.1	47.4	55.4	61.1
JÂA-Net* [101]	47.2	41.6	49.1	77.2	77.5	82.9	85.8	63.4	50.8	62.5	47.2	52.7	61.5
SRERL [51]	46.9	45.3	55.6	77.1	78.4	83.5	87.6	60.6	52.2	63.9	47.1	53.3	62.9
HMP-PS [105]	53.1	46.1	56.0	76.5	76.9	82.1	86.4	64.8	51.5	63.0	49.9	54.5	63.4
SEV-Net [130]	58.2	50.4	58.3	81.9	73.9	87.8	87.5	61.6	52.6	62.2	44.6	47.6	63.9
PIAP [113]	54.2	47.1	54.0	79.0	78.2	86.3	89.5	66.1	49.7	63.2	49.9	52.0	64.1
ATCM [42]	51.7	49.3	61.0	77.8	79.5	82.9	86.3	67.6	51.9	63.0	43.7	56.3	64.2
GraphAU (R-50) [79]	53.7	46.9	59.0	78.5	80.0	84.4	87.8	67.3	52.5	63.2	50.6	52.4	64.7
$ReCoT^*$ [61]	51.5	47.8	58.9	79.2	80.2	84.9	88.4	61.6	53.3	64.6	51.8	55.4	64.8
EmoLA (Ours.)	57.4	52.4	61.0	78.1	77.8	81.9	89.5	60.5	49.3	64.9	46.0	52.4	64.2

Table 6: F1 score (in %) results of 10 AUs on GFT. The results with * are taken from [108]. The others are taken directly from original papers. "ft" stands for finetune.

Method/AU	1	2	4	6	10	12	14	15	23	24	Avg.
EACNet [57]	15.5	56.6	0.1	81.0	76.1	84.0	0.1	38.5	57.8	51.2	46.1
TCAE [65]	43.9	49.5	6.3	71.0	76.2	79.5	10.7	28.5	34.5	41.7	44.2
ARL [102]	51.9	45.9	13.7	79.2	75.5	82.8	0.1	44.9	59.2	47.5	50.1
MoCo (ft)* [36]	45.3	48.2	20.3	80.7	78.8	78.1	22.6	46.0	53.9	50.3	52.4
Temporal Ranking (ft)* [77]	58.8	56.8	33.2	72.5	76.2	80.8	19.9	46.8	55.2	47.3	54.7
JÂANet [101]	46.5	49.3	19.2	79.0	75.0	84.8	44.1	33.5	54.9	50.7	53.7
EmoCo [108]	65.9	55.9	40.7	83.1	75.1	81.4	21.3	48.5	58.0	56.5	58.6
EmoLA (Ours.)	69.8	59.1	52.8	85.3	73.0	85.3	32.3	47.6	63.1	52.2	62.1

baselines mainly employ Vicuna [15] as LLM decoder and incorporate instruction tuning. We reproduce all the baselines using their open-sourced codes. For fair comparison, all the models are trained for 1 epoch and train on two tasks individually. Except for our EmoLA, all the other baselines are finetuned based on the pretrained MLLMs. More details about baselines can refer to the Appendix.

It is noteworthy that EmoLA achieves the best results in two tasks with fewer tuning parameters compared to other MLLMs. EmoLA finetunes merely 10% of the parameters compared to LLaVA-1.5, yet it achieves better performance on FABA-Bench, which can be attributed to two aspects. For one thing, EmoLA only finetunes LoRA which is much more efficient than tuning the entire LLM decoder. For another, EmoLA incorporates a facial prior expert to extract the facial structure knowledge. This feature compensates for the information overlooked by the visual encoder in FABA tasks. Moreover, we can also see that due to the strong capability of LLM, the language generation ability of Shikra, LLaVA-1.5 and EmoLA is comparable, making it unreasonable to evaluate these models merely based on the NLG metrics. Our metric, REGE, accounts for both

Table 7: Comparison on FABA-bench. All the baselines are reproduced by their open-sourced codes. S_{re} in emotion and AU tasks stands for accuracy and average F1 scores, respectively. S_{ge} and S_{rege} are ROUGE and our REGE score. All scores are in %.

Methods	E	Emotion		AU														
Methods	S_{re}	S_{ge}	S_{rege}	1	2	4	5	6	10	12	17	24	25	26	43	S_{re}	S_{ge}	S_{rege}
MiniGPT4-v2 [10]	58.2	19.6	77.8	47.9	35.5	42.3	32.7	29.2	6.6	10.3	0.0	2.5	0.1	0.0	0.0	17.9	19.9	37.8
mPLUG-Owl2 [131]	53.6	28.4	82.0	72.3	17.5	75.2	54.2	75.6	0.0	13.0	0.0	0.0	3.9	18.2	0.0	27.5	28.2	55.7
Shikra [11]	62.5	32.1	94.6	70.6	33.9	76.6	63.3	57.8	43.4	58.0	53.0	54.1	68.5	42.4	0.0	51.8	34.8	86.6
LLaVA-1.5 [70]	62.3	31.6	93.9	74.2	32.7	76.5	67.9	63.6	41.0	61.0	53.4	54.1	67.5	43.5	50.0	57.1	34.3	91.4
EmoLA (Ours.)	64.5	31.7	96.2	72.8	37.3	79.9	67.3	69.9	41.7	63.6	56.8	55.6	73.4	56.8	0.0	56.3	35.2	91.5

the recognition and generation abilities of MLLMs in FABA, offering a more comprehensive evaluation of their performance.

5.3 Ablation study

We explore the effectiveness of the facial landmark token in Tab. 8, and the tuning strategies in Tab. 9. More ablations can be found in Appendix, e.g., multitask training, the position of the prior token, etc. Fig. 8 shows the capability of EmoLA on FABA tasks, with additional examples available in the Appendix.

Effectiveness of prior token. As shown in Tab. 8, if we input only the facial prior token H_p to EmoLA without visual tokens, there will be a significant decline in recognition performance. However, thanks to the powerful capability of LLM decoder, the generalization ability does not suffer greatly. Furthermore, even with just a single token fed into the decoder, the model

Table 8: Influence of prior token. H_v and H_p indicate visual tokens and the facial prior token, respectively.

Tok	cens	Emo	otion	AU			
H_v	H_p	S_{re}	S_{ge}	S_{re}	S_{ge}		
	1	41.2	29.9	40.5	33.7		
✓		62.5	32.1	55.3	34.8		
✓	✓	64.5	31.7	56.3	35.2		

retains a certain level of recognition capability. This demonstrates that the land-mark prior contains representative features for FABA tasks.

Compared to merely inputting the prior token H_p , incorporating visual tokens H_v will enhance the recognition performance since these visual tokens carry more detailed face information. Incorporating both visual tokens and prior tokens leads to further improvements with eligible performance dropping in text generation. This suggests that the face prior token supplies complementary information of visual tokens. We suppose this additional information includes facial structure details introduced by landmark facial prior extractor.

Influence of tuning strategy. We also investigate the impact of tuning strategies on two modules requiring optimization, h_{ϕ} and g_{θ} , as shown in Tab. 9. When only the prior projector g_{θ} is fine-tuned, with all other parameters kept frozen, we observe that while the overall recognition performance is relatively low, it still surpasses that achieved by using only a single prior token. This suggests the effectiveness of tuning the prior projector g_{θ} , since the parameters of g_{θ} only take a small portion of the entire parameters.

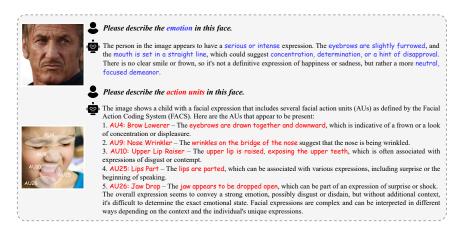


Fig. 8: An example of EmoLA's capability on FABA tasks.

We also tried exclusively tuning LoRA with visual projector $h_{\phi} + h_{\gamma}$ while keeping other parameters frozen. This is different from merely inputting the visual tokens H_v since it also includes a prior token H_p generated by the randomly initialized prior projector g_{θ} . The results indicate that tuning $h_{\phi} + h_{\gamma}$ significantly improves the performance by better aligning the output of the LLM. It can also be observed that the per-

Table 9: Influence of the tuning strategy. $h_{\phi} + h_{\gamma}$ and g_{θ} indicate the LoRA module with visual projector and the prior projector, respectively.

Module	Emo	otion	AU			
$h_{\phi} + h_{\gamma}$	g_{θ}	S_{re}	S_{ge}	S_{re}	S_{ge}	
	/	44.9	29.6	47.7	34.0	
✓		63.0	32.1	55.6	34.9	
✓	✓	64.5	31.7	56.3	35.2	

formance slightly exceeds that achieved using only the visual tokens, which also emphasizes the value of the prior token. Again, the efficacy is further augmented by simultaneously fine-tuning both $h_{\theta} + h_{\gamma}$ and the prior projector g_{θ} .

6 Conclusion

In this paper, to address the challenges raised in FABA tasks when employing MLLMs, we proposed an instruction-following FABA dataset FABA-Instruct by means of GPT4. Based on this dataset, we introduced a benchmark FABA-Bench to comprehensively evaluate the FABA models on instruction-following data. Furthermore, we presented an instruction-tuned MLLM EmoLA for FABA, which is efficient and effective by tuning LoRA on a pre-trained MLLM and incorporating a facial prior expert. Extensive experiments across four traditional FABA datasets and our FABA-Bench demonstrate the effectiveness of EmoLA. In the future, we intend to broaden our method to additional facial-related tasks, e.g., face detection, face generation, etc. Moreover, incorporating other facial prior features holds the potential for performance improvement. Expanding our EmoLA from 2D face images to video streams also presents a promising avenue for future research.

Acknowledgement: Yifan Li, Wentao Bao, and Yu Kong are partially supported by NSF Awards 1949694 and 2040209. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF. We also would like to express our deepest gratitude to Dr. Junwen Chen for her invaluable support and contribution to this research.

References

- Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a visual language model for few-shot learning. In: Adv. Neural Inform. Process. Syst. (2022) 4
- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al.: Qwen technical report. arXiv preprint arXiv:2309.16609 (2023)
- 3. Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A frontier large vision-language model with versatile abilities. arXiv preprint arXiv:2308.12966 (2023) 5
- Barrett, L.F., Mesquita, B., Ochsner, K.N., Gross, J.J.: The experience of emotion. Annu. Rev. Psychol. 58, 373–403 (2007)
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. vol. 33, pp. 1877–1901 (2020) 25
- Chang, Y., Wang, S.: Knowledge-driven self-supervised representation learning for facial action unit recognition. In: CVPR. pp. 20417–20426 (2022) 4
- 7. Chen, F., Han, M., Zhao, H., Zhang, Q., Shi, J., Xu, S., Xu, B.: X-llm: Bootstrapping advanced large language models by treating multi-modalities as foreign languages. arXiv preprint arXiv:2305.04160 (2023) 5
- 8. Chen, F., Shao, J., Zhu, S., Shen, H.T.: Multivariate, multi-frequency and multi-modal: Rethinking graph neural networks for emotion recognition in conversation. In: CVPR. pp. 10761–10770 (2023) 4
- 9. Chen, J., Zhang, A., Shi, X., Li, M., Smola, A., Yang, D.: Parameter-efficient fine-tuning design spaces. arXiv preprint arXiv:2301.01821 (2023) 5
- Chen, J., Zhu, D., Shen, X., Li, X., Liu, Z., Zhang, P., Krishnamoorthi, R., Chandra, V., Xiong, Y., Elhoseiny, M.: Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. arXiv preprint arXiv:2310.09478 (2023) 11, 13, 25, 27
- 11. Chen, K., Zhang, Z., Zeng, W., Zhang, R., Zhu, F., Zhao, R.: Shikra: Unleashing multimodal llm's referential dialogue magic. arXiv preprint arXiv:2306.15195 (2023) 2, 4, 11, 13, 25, 27
- 12. Chen, R., Zhang, H., Liang, S., Li, J., Cao, X.: Less is more: Fewer interpretable region via submodular subset selection. In: ICLR (2024) 7
- 13. Chen, T., Saxena, S., Li, L., Fleet, D.J., Hinton, G.: Pix2seq: A language modeling framework for object detection. In: ICLR (2021) 2
- 14. Chen, Z., Wu, J., Wang, W., Su, W., Chen, G., Xing, S., Muyan, Z., Zhang, Q., Zhu, X., Lu, L., et al.: Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. arXiv preprint arXiv:2312.14238 (2023) 5

- Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., Stoica, I., Xing, E.P.: Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality (March 2023), https://lmsys.org/ blog/2023-03-30-vicuna/ 4, 10, 12
- Cho, J., Lei, J., Tan, H., Bansal, M.: Unifying vision-and-language tasks via text generation. In: ICML. pp. 1931–1942 (2021)
- 17. Chung, H.W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., et al.: Scaling instruction-finetuned language models. arXiv preprint arXiv:2210.11416 (2022) 4, 24
- 18. Corneanu, C., Madadi, M., Escalera, S.: Deep structure inference network for facial action unit recognition. In: ECCV. pp. 298–313 (2018) 11, 12
- Cui, Z., Kuang, C., Gao, T., Talamadupula, K., Ji, Q.: Biomechanics-guided facial action unit detection through force modeling. In: CVPR. pp. 8694–8703 (2023) 4
- Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning (2023) 5
- Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: CVPR. pp. 4690–4699 (2019)
- 22. Dettmers, T., Pagnoni, A., Holtzman, A., Zettlemoyer, L.: Qlora: Efficient fine-tuning of quantized llms. Adv. Neural Inform. Process. Syst. **36** (2024) 5
- Dhall, A., Goecke, R., Lucey, S., Gedeon, T.: Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark. In: IEEE Int. Conf. Comput. Vis. Worksh. pp. 2106–2112 (2011) 6
- Dong, R., Han, C., Peng, Y., Qi, Z., Ge, Z., Yang, J., Zhao, L., Sun, J., Zhou, H., Wei, H., et al.: Dreamllm: Synergistic multimodal comprehension and creation. arXiv preprint arXiv:2309.11499 (2023) 4
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2020) 4
- Ekman, P., Friesen, W.V.: Facial action coding system. Environmental Psychology & Nonverbal Behavior (1978) 1, 4, 7
- 27. Ekman, P., et al.: Basic emotions. Handbook of cognition and emotion $\bf 98(45-60)$, $\bf 16\ (1999)\ 4$
- 28. Fabian Benitez-Quiroz, C., Srinivasan, R., Martinez, A.M.: Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In: CVPR. pp. 5562–5570 (2016) 3, 5, 7, 9
- Gatt, A., Krahmer, E.: Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. Jour. Art. Intel. Resea. 61, 65–170 (2018) 8
- 30. Girard, J.M., Chu, W.S., Jeni, L.A., Cohn, J.F.: Sayette group formation task (gft) spontaneous facial expression database. In: IEEE FG. pp. 581–588 (2017) 5, 6, 7, 10, 26
- 31. Goodfellow, I.J., Erhan, D., Carrier, P.L., Courville, A., Mirza, M., Hamner, B., Cukierski, W., Tang, Y., Thaler, D., Lee, D.H., et al.: Challenges in representation learning: A report on three machine learning contests. In: Adv. Neural Inform. Process. Syst. pp. 117–124 (2013) 5, 6
- 32. Grandjean, D., Sander, D., Scherer, K.R.: Conscious emotional experience emerges as a function of multilevel, appraisal-driven response synchronization. Consciousness and cognition 17(2), 484–495 (2008) 2, 4
- 33. Guo, D., Rush, A.M., Kim, Y.: Parameter-efficient transfer learning with diff pruning. arXiv preprint arXiv:2012.07463 (2020) 5

- 34. Haidt, J., Keltner, D.: Culture and facial expression: Open-ended methods find more expressions and a gradient of recognition. Cognition & Emotion 13(3), 225–266 (1999) 7
- He, J., Zhou, C., Ma, X., Berg-Kirkpatrick, T., Neubig, G.: Towards a unified view of parameter-efficient transfer learning. arXiv preprint arXiv:2110.04366 (2021)
- 36. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: CVPR. pp. 9729–9738 (2020) 12
- 37. He, S., Ding, L., Dong, D., Zhang, M., Tao, D.: Sparseadapter: An easy approach for improving the parameter-efficiency of adapters. arXiv preprint arXiv:2210.04284 (2022) 5
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., Gelly, S.: Parameter-efficient transfer learning for nlp. In: ICML. pp. 2790–2799 (2019) 5
- Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021) 3, 5
- Iyer, S., Lin, X.V., Pasunuru, R., Mihaylov, T., Simig, D., Yu, P., Shuster, K., Wang, T., Liu, Q., Koura, P.S., et al.: Opt-iml: Scaling language model instruction meta-learning through the lens of generalization. arXiv preprint arXiv:2212.12017 (2022) 24
- 41. Izard, C.E.: Human emotions. Springer Science & Business Media (2013) 4
- 42. Jacob, G.M., Stenger, B.: Facial action unit detection with transformers. In: CVPR. pp. 7680–7689 (2021) 4, 11, 12
- 43. Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., Casas, D.d.l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al.: Mistral 7b. arXiv preprint arXiv:2310.06825 (2023) 2, 4
- 44. Jiang, X., Zong, Y., Zheng, W., Tang, C., Xia, W., Lu, C., Liu, J.: Dfew: A large-scale database for recognizing dynamic facial expressions in the wild. In: ACM MM. pp. 2881–2889 (2020) 5, 6
- Karimi Mahabadi, R., Henderson, J., Ruder, S.: Compacter: Efficient low-rank hypercomplex adapter layers. Adv. Neural Inform. Process. Syst. 34, 1022–1035 (2021) 5
- 46. Kollias, D., Schulc, A., Hajiyev, E., Zafeiriou, S.: Analysing affective behavior in the first abaw 2020 competition. In: IEEE FG. pp. 637–643 (2020) 1
- 47. Kollias, D., Zafeiriou, S.: Aff-wild2: Extending the aff-wild database for affect recognition. arXiv preprint arXiv:1811.07770 (2018) 3, 5
- 48. Lai, X., Tian, Z., Chen, Y., Li, Y., Yuan, Y., Liu, S., Jia, J.: Lisa: Reasoning segmentation via large language model. arXiv preprint arXiv:2308.00692 (2023)
- 49. Lester, B., Al-Rfou, R., Constant, N.: The power of scale for parameter-efficient prompt tuning. arXiv preprint arXiv:2104.08691 (2021) 5
- Li, B., Zhang, Y., Chen, L., Wang, J., Yang, J., Liu, Z.: Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726 (2023)
- Li, G., Zhu, X., Zeng, Y., Wang, Q., Lin, L.: Semantic relationships guided representation learning for facial action unit recognition. In: AAAI. pp. 8594–8601 (2019) 4, 11, 12
- 52. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) 2, 5

- Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: ICML. pp. 12888– 12900 (2022)
- 54. Li, S., Deng, W.: Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. IEEE TIP **28**(1), 356–370 (2019) **3**
- Li, S., Deng, W.: Deep facial expression recognition: A survey. IEEE Trans. Affect. Comput. 13(3), 1195–1215 (2020) 1
- Li, S., Deng, W., Du, J.: Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild. In: CVPR. pp. 2584–2593 (2017) 3, 5, 6, 10, 26
- 57. Li, W., Abtahi, F., Zhu, Z., Yin, L.: Eac-net: Deep nets with enhancing and cropping for facial action unit detection. IEEE TPAMI **40**(11), 2583–2596 (2018) 11, 12
- 58. Li, X.L., Liang, P.: Prefix-tuning: Optimizing continuous prompts for generation. arXiv preprint arXiv:2101.00190 (2021) 5
- Li, X., Behpour, S., Doan, T.L., He, W., Gou, L., Ren, L.: Up-dp: Unsupervised prompt learning for data pre-selection with vision-language models. In: Adv. Neural Inform. Process. Syst. vol. 36 (2024) 5
- 60. Li, X., Pan, D., Li, C., Qiang, Y., Zhu, D.: Negative flux aggregation to estimate feature attributions. In: IJCAI (2023) 7
- Li, Y., Han, H., Shan, S., Ji, Z., Bai, J., Chen, X.: Recot: Regularized co-training for facial action unit recognition with noisy labels. In: BMVC (2023) 4, 10, 11, 12
- 62. Li, Y., Sun, H., Liu, Z., Han, H., Shan, S.: Affective behaviour analysis using pretrained model with facial prior. In: Eur. Conf. Comput. Vis. Worksh. pp. 19–30 (2022) 4
- 63. Li, Y., Wang, Y., Cui, Z.: Decoupled multimodal distilling for emotion recognition. In: CVPR. pp. 6631–6640 (2023) 4
- 64. Li, Y., Zeng, J., Shan, S., Chen, X.: Occlusion aware facial expression recognition using cnn with attention mechanism. IEEE TIP **28**(5), 2439–2450 (2018) **2**
- Li, Y., Zeng, J., Shan, S., Chen, X.: Self-supervised representation learning from videos for facial action unit detection. In: CVPR. pp. 10924–10933 (2019) 4, 12
- 66. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out. pp. 74-81 (2004) 8
- 67. Lin, J., Yin, H., Ping, W., Lu, Y., Molchanov, P., Tao, A., Mao, H., Kautz, J., Shoeybi, M., Han, S.: Vila: On pre-training for visual language models. arXiv preprint arXiv:2312.07533 (2023) 4
- 68. Lin, Z., Madotto, A., Fung, P.: Exploring versatile generative language model via parameter-efficient transfer learning. arXiv preprint arXiv:2004.03829 (2020) 5
- Liu, H., Tam, D., Muqeeth, M., Mohta, J., Huang, T., Bansal, M., Raffel, C.A.: Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. Adv. Neural Inform. Process. Syst. 35, 1950–1965 (2022) 5
- 70. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744 (2023) 2, 3, 5, 10, 11, 13, 27
- 71. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. Adv. Neural Inform. Process. Syst. **36** (2024) **2**, **5**, **25**
- Liu, M., Shan, S., Wang, R., Chen, X.: Learning expressionlets on spatio-temporal manifold for dynamic facial expression recognition. In: CVPR. pp. 1749–1756 (2014) 7

- 73. Liu, Q., Wu, X., Zhao, X., Zhu, Y., Xu, D., Tian, F., Zheng, Y.: Moelora: An moe-based parameter efficient fine-tuning method for multi-task medical applications. arXiv preprint arXiv:2310.18339 (2023) 5
- 74. Liu, Y., Dai, W., Feng, C., Wang, W., Yin, G., Zeng, J., Shan, S.: Mafw: A large-scale, multi-modal, compound affective database for dynamic facial expression recognition in the wild. In: ACM MM. pp. 24–32 (2022) 3, 5, 6
- 75. Lopes, A.T., De Aguiar, E., De Souza, A.F., Oliveira-Santos, T.: Facial expression recognition with convolutional neural networks: coping with few data and the training sample order. PR **61**, 610–628 (2017) 1
- Lu, H., Niu, X., Wang, J., Wang, Y., Hu, Q., Tang, J., Zhang, Y., Yuan, K., Huang, B., Yu, Z., et al.: Gpt as psychologist? preliminary evaluations for gpt-4v on visual affective computing. arXiv preprint arXiv:2403.05916 (2024) 6
- 77. Lu, L., Tavabi, L., Soleymani, M.: Self-supervised learning for facial action unit recognition through temporal consistency. In: BMVC (2020) 12
- 78. Lucey, P., Cohn, J.F., Kanade, T., Saragih, J., Ambadar, Z., Matthews, I.: The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In: CVPRW. pp. 94–101 (2010) 5, 6
- Luo, C., Song, S., Xie, W., Shen, L., Gunes, H.: Learning multi-dimensional edge feature-based au relation graph for facial action unit recognition. arXiv preprint arXiv:2205.01782 (2022) 11, 12
- 80. Mahabadi, R.K., Ruder, S., Dehghani, M., Henderson, J.: Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. arXiv preprint arXiv:2106.04489 (2021) 5
- 81. Mao, J., Xu, R., Yin, X., Chang, Y., Nie, B., Huang, A.: Poster v2: A simpler and stronger facial expression recognition network. arXiv preprint arXiv:2301.12149 (2023) 10
- 82. Mao, Y., Mathias, L., Hou, R., Almahairi, A., Ma, H., Han, J., Yih, W.t., Khabsa, M.: Unipelt: A unified framework for parameter-efficient language model tuning. arXiv preprint arXiv:2110.07577 (2021) 5
- 83. Martinez, B., Valstar, M.F., Jiang, B., Pantic, M.: Automatic analysis of facial actions: A survey. IEEE Trans. Affect. Comput. 10(3), 325–347 (2017) 1
- 84. Mavadati, S.M., Mahoor, M.H., Bartlett, K., Trinh, P., Cohn, J.F.: Disfa: A spontaneous facial action intensity database. IEEE Trans. Affect. Comput. 4(2), 151–160 (2013) 3, 5, 6, 7, 10, 26
- 85. Mollahosseini, A., Hasani, B., Mahoor, M.H.: Affectnet: A database for facial expression, valence, and arousal computing in the wild. IEEE Trans. Affect. Comput. **10**(1), 18–31 (2017) **4**, **5**, **6**, **26**
- 86. Niu, X., Han, H., Shan, S., Chen, X.: Multi-label co-regularization for semi-supervised facial action unit recognition. In: Adv. Neural Inform. Process. Syst. pp. 909–919 (2019) 4
- 87. Niu, X., Han, H., Yang, S., Huang, Y., Shan, S.: Local relationship learning with person-specific shape regularization for facial action unit detection. In: CVPR. pp. 11917–11926 (2019) 2, 4, 10, 11
- 88. OpenAI, R.: Gpt-4 technical report. arxiv 2303.08774. View in Article $\bf 2$, 13 (2023) $\bf 24$
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. In: Adv. Neural Inform. Process. Syst. vol. 35, pp. 27730–27744 (2022) 24
- 90. Pantic, M., Valstar, M., Rademaker, R., Maat, L.: Web-based database for facial expression analysis. In: ICME. pp. 5–pp (2005) 5, 6

- 91. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Asso. Comput. Ling. pp. 311–318 (2002) 8
- 92. Pfeiffer, J., Kamath, A., Rücklé, A., Cho, K., Gurevych, I.: Adapterfusion: Non-destructive task composition for transfer learning. arXiv preprint arXiv:2005.00247 (2020) 5
- 93. Pixton, T.S.: Happy to see me, aren't you, sally? signal detection analysis of emotion detection in briefly presented male and female faces. Scandinavian Journal of Psychology **52**(4), 361–368 (2011) 2
- 94. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML. pp. 8748–8763 (2021) 3, 4, 9
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li,
 W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. JMLR 21(140), 1–67 (2020) 25
- Rücklé, A., Geigle, G., Glockner, M., Beck, T., Pfeiffer, J., Reimers, N., Gurevych,
 I.: Adapterdrop: On the efficiency of adapters in transformers. arXiv preprint arXiv:2010.11918 (2020) 5
- 97. Russell, J.A.: A circumplex model of affect. Journal of personality and social psychology 39(6), 1161 (1980) 4
- 98. Saneiro, M., Santos, O.C., Salmeron-Majadas, S., Boticario, J.G., et al.: Towards emotion detection in educational scenarios from facial expressions and body movements through multimodal approaches. The Scientific World Journal **2014** (2014) 2
- 99. Sankaran, N., Mohan, D.D., Setlur, S., Govindaraju, V., Fedorishin, D.: Representation learning through cross-modality supervision. In: IEEE FG. pp. 1–8 (2019) 11, 12
- 100. Savchenko, A.V.: Frame-level prediction of facial expressions, valence, arousal and action units for mobile devices. arXiv preprint arXiv:2203.13436 (2022) 4
- 101. Shao, Z., Liu, Z., Cai, J., Ma, L.: Jaa-net: Joint facial action unit detection and face alignment via adaptive attention. IJCV 129(2), 321–340 (2021) 2, 10, 11, 12
- 102. Shao, Z., Liu, Z., Cai, J., Wu, Y., Ma, L.: Facial action unit detection using attention and relation learning. IEEE Trans. Affect. Comput. (2019) 4, 11, 12
- 103. Shen, J., Wang, H., Gui, S., Tan, J., Wang, Z., Liu, J.: {UMEC}: Unified model and embedding compression for efficient recommendation systems. In: International Conference on Learning Representations (2021), https://openreview.net/forum?id=BM---bH_RSh 5
- 104. Song, T., Cui, Z., Wang, Y., Zheng, W., Ji, Q.: Dynamic probabilistic graph convolution for facial action unit intensity estimation. In: CVPR. pp. 4845–4854 (2021) 4
- Song, T., Cui, Z., Zheng, W., Ji, Q.: Hybrid message passing with performancedriven structures for facial action unit detection. In: CVPR. pp. 6267–6276 (2021) 11, 12
- Sun, B., Cao, S., Li, D., He, J., Yu, L.: Dynamic micro-expression recognition using knowledge distillation. IEEE Trans. Affect. Comput. 13(2), 1037–1043 (2020)
- 107. Sun, Q., Yu, Q., Cui, Y., Zhang, F., Zhang, X., Wang, Y., Gao, H., Liu, J., Huang, T., Wang, X.: Generative pretraining in multimodality. arXiv preprint arXiv:2307.05222 (2023) 5
- 108. Sun, X., Zeng, J., Shan, S.: Emotion-aware contrastive learning for facial action unit detection. In: FG. pp. 01–08 (2021) 4, 11, 12

- 109. Sung, Y.L., Nair, V., Raffel, C.A.: Training neural networks with fixed sparse masks. Adv. Neural Inform. Process. Syst. 34, 24193–24205 (2021) 5
- 110. Tan, Z., Beigi, A., Wang, S., Guo, R., Bhattacharjee, A., Jiang, B., Karami, M., Li, J., Cheng, L., Liu, H.: Large language models for data annotation: A survey. arXiv preprint arXiv:2402.13446 (2024) 6
- Tan, Z., Chen, T., Zhang, Z., Liu, H.: Sparsity-guided holistic explanation for llms with interpretable inference-time intervention. In: AAAI. pp. 21619–21627 (2024) 24
- 112. Tan, Z., Cheng, L., Wang, S., Bo, Y., Li, J., Liu, H.: Interpreting pretrained language models via concept bottlenecks. arXiv preprint arXiv:2311.05014 (2023) 24
- 113. Tang, Y., Zeng, W., Zhao, D., Zhang, H.: Piap-df: Pixel-interested and anti person-specific facial action unit detection net with discrete feedback learning. In: ICCV. pp. 12899–12908 (2021) 4, 11, 12
- 114. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023) 2, 4
- 115. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023) 2
- Wang, K., Peng, X., Yang, J., Meng, D., Qiao, Y.: Region attention networks for pose and occlusion robust facial expression recognition. IEEE TIP 29, 4057–4069 (2020) 2, 4, 10
- 117. Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., et al.: Cogvlm: Visual expert for pretrained language models. arXiv preprint arXiv:2311.03079 (2023) 4
- 118. Wang, Y., Sun, Y., Huang, Y., Liu, Z., Gao, S., Zhang, W., Ge, W., Zhang, W.: Ferv39k: A large-scale multi-scene dataset for facial expression recognition in videos. In: CVPR. pp. 20922–20931 (2022) 5, 6
- Wang, Y., Agarwal, S., Mukherjee, S., Liu, X., Gao, J., Awadallah, A.H., Gao, J.: Adamix: Mixture-of-adaptations for parameter-efficient model tuning. arXiv preprint arXiv:2210.17451 (2022) 5
- 120. Wang, Y., Chen, W., Han, X., Lin, X., Zhao, H., Liu, Y., Zhai, B., Yuan, J., You, Q., Yang, H.: Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. arXiv preprint arXiv:2401.06805 (2024) 2
- 121. Wang, Z., Zeng, F., Liu, S., Zeng, B.: Oaenet: Oriented attention ensemble for accurate facial expression recognition. PR 112, 107694 (2021) 4
- 122. Wen, Z., Lin, W., Wang, T., Xu, G.: Distract your attention: Multi-head cross attention network for facial expression recognition. Biomimetics 8(2), 199 (2023) 4, 10
- 123. Wu, S., Fei, H., Qu, L., Ji, W., Chua, T.S.: NExT-GPT: Any-to-any multimodal llm. CoRR abs/2309.05519 (2023) 4
- 124. Xie, Q., Luong, M.T., Hovy, E., Le, Q.V.: Self-training with noisy student improves imagenet classification. In: CVPR. pp. 10687–10698 (2020) 4
- 125. Xie, S., Hu, H., Wu, Y.: Deep multi-path convolutional neural network joint with salient region attention for facial expression recognition. PR 92, 177–191 (2019)
- 126. Xu, L., Xie, H., Qin, S.Z.J., Tao, X., Wang, F.L.: Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. arXiv preprint arXiv:2312.12148 (2023) 5

- 127. Xu, R., Luo, F., Zhang, Z., Tan, C., Chang, B., Huang, S., Huang, F.: Raise a child in large language model: Towards effective and generalizable fine-tuning. arXiv preprint arXiv:2109.05687 (2021) 5
- 128. Xue, F., Wang, Q., Tan, Z., Ma, Z., Guo, G.: Vision transformer with attentive pooling for robust facial expression recognition. IEEE Trans. Affect. Comput. (2022) 4, 10
- 129. Yan, W.J., Li, X., Wang, S.J., Zhao, G., Liu, Y.J., Chen, Y.H., Fu, X.: Casme ii: An improved spontaneous micro-expression database and the baseline evaluation. PloS one **9**(1), e86041 (2014) 5
- Yang, H., Yin, L., Zhou, Y., Gu, J.: Exploiting semantic embedding and visual feature for facial action unit detection. In: CVPR. pp. 10482–10491 (2021) 4, 11, 12
- 131. Ye, Q., Xu, H., Ye, J., Yan, M., Liu, H., Qian, Q., Zhang, J., Huang, F., Zhou, J.: mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. arXiv preprint arXiv:2311.04257 (2023) 11, 13, 25, 27
- 132. Yin, Y., Chang, D., Song, G., Sang, S., Zhi, T., Liu, J., Luo, L., Soleymani, M.: Fg-net: Facial action unit detection with generalizable pyramidal features. In: Wint. Appl. Comput. Vis. pp. 6099–6108 (2024) 4
- 133. You, H., Zhang, H., Gan, Z., Du, X., Zhang, B., Wang, Z., Cao, L., Chang, S.F., Yang, Y.: Ferret: Refer and ground anything anywhere at any granularity. In: ICLR (2023) 2
- 134. Yu, S., Chen, T., Shen, J., Yuan, H., Tan, J., Yang, S., Liu, J., Wang, Z.: Unified visual transformer compression. arXiv preprint arXiv:2203.08243 (2022) 5
- 135. Zaken, E.B., Ravfogel, S., Goldberg, Y.: Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. arXiv preprint arXiv:2106.10199 (2021) 5
- 136. Zhang, D., Yu, Y., Li, C., Dong, J., Su, D., Chu, C., Yu, D.: Mm-llms: Recent advances in multimodal large language models. In: arXiv preprint arXiv:2401.13601 (2024) 2
- 137. Zhang, Q., Chen, M., Bukharin, A., He, P., Cheng, Y., Chen, W., Zhao, T.: Adaptive budget allocation for parameter-efficient fine-tuning. arXiv preprint arXiv:2303.10512 (2023) 5
- 138. Zhang, S., Dong, L., Li, X., Zhang, S., Sun, X., Wang, S., Li, J., Hu, R., Zhang, T., Wu, F., et al.: Instruction tuning for large language models: A survey. arXiv preprint arXiv:2308.10792 (2023) 24
- 139. Zhang, S., Pan, Y., Wang, J.Z.: Learning emotion representations from verbal and nonverbal communication. In: CVPR. pp. 18993–19004 (2023) 4
- 140. Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X.V., et al.: Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 (2022) 4, 25
- Zhang, X., Yang, H., Wang, T., Li, X., Yin, L.: Multimodal channel-mixing: Channel and spatial masked autoencoder on facial action unit detection. In: Wint. Appl. Comput. Vis. pp. 6077–6086 (2024) 4
- 142. Zhang, X., Yin, L., Cohn, J.F., Canavan, S., Reale, M., Horowitz, A., Liu, P.: A high-resolution spontaneous 3d dynamic facial expression database. In: IEEE FG. pp. 1–6 (2013) 3
- 143. Zhang, X., Yin, L., Cohn, J.F., Canavan, S., Reale, M., Horowitz, A., Liu, P., Girard, J.M.: Bp4d-spontaneous: A high-resolution spontaneous 3d dynamic facial expression database. Ima. Vis. Comput. 32(10), 692–706 (2014) 3, 5, 6, 7, 9, 10, 26

- Zhang, Y., Wang, C., Deng, W.: Relative uncertainty learning for facial expression recognition. In: Adv. Neural Inform. Process. Syst. pp. 17616–17627 (2021) 10
- Zhang, Y., Wang, C., Ling, X., Deng, W.: Learn from all: Erasing attention consistency for noisy label facial expression recognition. In: ECCV. pp. 418–434 (2022)
 4, 10
- 146. Zhang, Z., Wang, L., Yang, J.: Weakly supervised video emotion detection and prediction via cross-modal temporal erasing network. In: CVPR. pp. 18888–18897 (2023) 4
- 147. Zhao, K., Chu, W.S., Martinez, A.M.: Learning facial action units from web images with scalable weakly supervised clustering. In: CVPR. pp. 2090–2099 (2018)
- 148. Zhao, K., Chu, W.S., Zhang, H.: Deep region and multi-label learning for facial action unit detection. In: CVPR. pp. 3391–3399 (2016) 11, 12, 26
- 149. Zhao, S., Li, Y., Yao, X., Nie, W., Xu, P., Yang, J., Keutzer, K.: Emotion-based end-to-end matching between image and music in valence-arousal space. In: ACM MM. pp. 2945–2954 (2020) 4
- 150. Zhao, Z., Liu, Q., Wang, S.: Learning deep global multi-scale and local attention features for facial expression recognition in the wild. IEEE TIP 30, 6544–6556 (2021) 10
- Zhao, Z., Liu, Q., Zhou, F.: Robust lightweight facial expression recognition network with label distribution training. In: AAAI. pp. 3510–3519 (2021) 10
- 152. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023) 2, 5, 25

Facial Affective Behavior Analysis with Instruction Tuning

Supplementary Material

This is the supplementary material of Facial Affective Behavior Analysis with Instruction Tuning.

A Limitations and negative effects

A.1 Limitations

This work also has its limitations. Firstly, we haven't tried other face-prior feature extractors except for landmark features. We leave this as an exploring research direction for future work. Secondly, our annotations from the training set also include noise induced by GPT-4V's hallucinations [111, 112] (see Appendix C.3), which may introduce some bias for models. We think our benchmark can be regarded as the learning from noisy descriptions task. Thirdly, the metric we proposed for evaluating the expression recognition ability may not reflect the nuanced expressions, which could be further improved.

A.2 Negative effects

There also exist some potential negative effects of our EmoLA. Firstly, privacy issues. Facial affective behavior analysis (FABA) may infringe on users' privacy, especially when deploying EmoLA on public spaces or systems. Individuals' faces may be unconsciously captured and recorded, leading to privacy concerns. Secondly, misdirections and misjudgments. Our EmoLA is not entirely accurate and may produce misjudged emotions. This can lead to misunderstanding or misdirection issues, especially in applications like security or judicial systems. Thirdly, technological abuse. Our method may be abused to suppress dissent or monitor political activities, thereby leading to social hazard or freedom restrictions.

B Background of instruction tuning

According to [138], instruction tuning refers to "the process of further training LLMs on a dataset consisting of (instruction, output) pairs in a supervised fashion, which bridges the gap between the next-word prediction objective of LLMs and the users' objective of having LLMs adhere to human instructions." The LLMs are typically trained on large general language corpora, which may differ from the users' objectives. As a result, to make the LLMs follow users' instructions, instruction tuning is proposed. For instance, InstructGPT [89]/ Chat-GPT [88], FLAN-T5 [17], OPT-IML [40] are tuned with instruction-following

data to enable their counterparts GPT-3 [5], T5 [95], OPT [140] have better generalization and few-shot abilities. Inspired by the success of instruction tuning for LLMs, LLaVA [71] attempts to extend this technique to the multimodal space, by introducing an MLP connector to map the visual tokens to language token space. Following that, other methods [10,11,131,152] also adopt this mechanism on multiple downstream tasks and achieve remarkable results.

C Annotation details

C.1 Instructions in FABA-Instruct

We tried different instructions in our FABA-Instruct datasets. Specifically, we adopt 100 carefully designed instructions for emotion and action unit (AU) recognition tasks, respectively. Some of these instructions for emotion and AU are delineated in Fig. 10a and Fig. 10b, respectively. As shown in these examples, these instructions are all with natural language format.

C.2 AU types in FABA-Instruct

As mentioned in the paper, we select 12 AUs of FABA-Instruct for evaluation, *i.e.*, AU1, AU2, AU4, AU5, AU6, AU10, AU12, AU17, AU24, AU25, AU26, AU43. Also, there exists in total of 19AUs in the training annotations, and the meaning of these AUs are given in Tab. 10.

C.3 The accuracy estimation of training annotations

There also exists some noise in the training set due to the hallucinations in GPT-4V. To estimate the label accuracy in training annotations, we randomly sample 200 samples from each task in FABA-Instruct, and manually re-annotate these samples. After that, we can roughly estimate the accuracy or F1 score of training annotations according to these manual annotations. Specifically, for the emotion task, we calculate accuracy by classifying the text into 7 classes, which has been introduced in the main content. For the AU task, we evaluate all the AUs using the F1 score.

For emotion annotations in FABA-Instruct, as shown in Tab. 11, the accuracy of training annotations is about 91%. For AU annotations, the average F1 is 76.1% for all the AUs (see Tab. 10). From the estimation results, it can be observed that although there are some noisy labels in both two tasks, the recognition performance of GPT-4V on two tasks is still high. Therefore, it's reasonable to use these annotations for further research. Our FABA-Bench can not only be utilized for FABA tasks but also be regarded as the learning from noisy annotations task. Some examples are be found in Fig. 9.

Table 10: The meaning of AUs in FABA-Instruct.

	N. / *
AUs	Meaning
AU1	inner brow raiser
AU2	outer brow raiser
AU4	brow lowerer
AU5	upper lid raiser
AU6	cheek raiser
AU7	lid tightener
AU9	nose wrinkler
AU10	upper lip raiser
AU12	lip corner puller
AU14	$\operatorname{dimpler}$
AU15	lip corner depressor
AU17	chin raiser
AU20	lip stretcher
AU23	lip tightener
AU24	lip pressor
AU25	lips part
AU26	jaw drop
AU27	mouth stretch
AU43	eyes closed

D Experiments

D.1 Details about the traditional FABA datasets

We utilize three traditional AU datasets (BP4D, DISFA, GFT) and one traditional emotion datasets (RAF-DB) for better evaluating the performance of our EmoLA compared with the previous SOTA methods. Furthermore, we also sampled around 20,000 face images from AffectNet to construct our FABA-Instruct dataset.

Emotion datasets. AffectNet [85] is a large-scale in-the-wild emotion database with more than 1M face images crawled from internet. We randomly sampled 20,000 images from this database for dataset construction. RAF-DB [56] is an facial expression dataset with 29,672 in-the-wild face images, which is annotated by 40 annotators with 7 single-label and two-tab compond emotion categories. In this paper, we mainly compare with other methods on the single-label subset. AU datasets. BP4D [143] is a spontaneous facial AU dataset with 328 videos from 41 subjects (23 females and 18 males). There are in total of 140,000 frames with 12 AUs (1, 2, 4, 6, 7, 10, 12, 14, 15, 17, 23, 24). DISFA [84] consists of 26 subjects (12 females and 14 males) and 130,000 frames with the AU intensities (from 0 to 5) annotations. Following [148], an AU with the intensity equal or greater than 2 is considered to be activated. 8 (1, 2, 4, 6, 9, 12, 25, 26) of the 12 AUs are utilized for evaluation. GFT [30] contains 96 participants from 32 three-person groups. 10 AUs (1, 2, 4, 6, 10, 12, 14, 15, 23, 24) are selected



The person in the image appears to be engaging in a conversation or reacting to something; their eyes are open wide, and the eyebrows are raised, which might suggest surprise or intense interest. However, without more context, it's challenging to precisely pinpoint the emotion as facial cues can be interpreted differently depending on the situation.



The individual in the image appears to have a somewhat furrowed brow and a tightness around the mouth, which might suggest feelings of anger, frustration, or intensity. There's also a possibility that the person is in the midst of speaking or reacting to something, as people's facial expressions can change rapidly depending on the context. Without additional context, it is hard to determine the exact emotion with certainty.



In the image, the facial expression appears to be a mix of emotions; it's not a clear-crut display of a single emotion. The person seems to be giving a somewhat intense gaze, which could be interpreted as focused, determined, or possibly pensive. The lips are slightly parted, which might suggest a moment of contemplation or surprise. However, without additional context, it's not definitive. Facial expressions can be complex and open to interpretation.



The person in the image appears to be showing an emotion that could be interpreted as **shock**, **surprise**, or **fear**. This is suggested by the wide-open eyes and mouth agape, combined with the slightly raised eyebrows. Such facial expressions often indicate a strong emotional reaction to something unexpected or alarming.

(a) Examples of wrong emotion annotations.



GT: AU4, AU7, AU23

The facial action coding system (FACS) is a tool for categorizing human facial movements by their appearance on the face. Based on the image provided, I can identify several facial action units (AUs) that might be in use by the person in the photo: 1. AU4: Brow Lowerer - The person's eyebrows appear to be drawn together and slightly downwards, which could be indicative of this action unit.

- AU9: Nose Wrinkler There seems to be a slight wrinkling of the nose, suggesting this action unit might be in play.
 AU15: Lip Corner Depressor The corners of the mouth appear to be pulled downward, which is consistent with this AU.
- 4. AU23: Lip Tightener The lips seem to be tight and pursed, which could hint at the activation of this action unit.
- 5. AU24: Lip Pressor If the lips are indeed pressed tightly against each other, this action unit could also be in use.
 These AUs may combine to convey emotions such as confusion, skepticism, or mild disapproval. Please note that deducing the exact emotional state or intent solely from facial cues without additional context might not be accurate, as facial expressions can be complex and subject to individual interpretation.



GT: AU1, AU5, AU25

In the image, the person seems to display a combination of facial expressions. The Facial Action Coding System (FACS) can be applied to describe the individual action units (AUs):

- AU1+AU2: Raised eyebrows, indicated by the visible lines across the forehead and the position of the eyebrows above the neutral position.
- AU4: Brow Lowerer, with the eyebrows drawn together, which often indicates confusion or concentration
 AU5: Upper eyelid raise, though this is not very prominent here.
- Aldís Cheek Raiser, which is associated with the orbicularis oculi muscle contracting and might be part of a squint.
 AU7: Lid Tightener may also be present to some degree, suggested by the tension around the eyes.
- AU12: Lip Corner Puller, as indicated by the corners of the lips being drawn back slightly, usually associated with smiling or a positive emotion. However, the intensity seems low as the expression does not look like a strong smile.
- AU25+AU26: Lips part slightly, which could imply a speech beginning or a slight surprise.
 Please note that without the complete context of the facial expression and additional movement data, this analysis is speculative based on static visual cues. Facial expressions are dynamic, and more accurate interpretations often require motion and context.

(b) Examples of wrong AU annotations.

Fig. 9: Wrong annotation examples on FABA-Instruct.

considering the challenges brought by head motion and occlusion. There are in total of 108000 training images and 24600 evaluating images in GFT.

D.2 Baseline details of FABA-Bench

As presented in the Experiment section, we reproduced four baselines which have similar arhitectures as our EmoLA on FABA-Bench, *i.e.*, MiniGPT4-V2 [10], mPLUG-Owl2 [131], Shikra [11], and LLaVA-1.5 [70]. LLaVA utilize a image encoder to obtain the visual tokens, and map the visual tokens to the language space through a linear layer. After that, the image tokens and language tokens are passed to an LLM decoder to generate descriptions. LLaVA undergoes a two-stage training process: initially, it exclusively trains the projector, followed

Table 11: Estimated accuracy of Emotion annotations in FABA-Instruct.

Total	Correct	Acc.
200	182	91.0

Table 12: Estimated F1 score of AU annotations in FABA-Instruct.

AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU12	AU15
92.2	95.6	84.2	90.7	82.6	65.9	61.3	58.1	84.3	70.6
AU17	AU2	0 AU	23 A	U24 .	AU25	AU26	AU27	AU43	Avg.

Table 13: The multitask performance of EmoLA on FABA-Bench.

Methods			on				
Wethods	S_{re}	S_{ge}	S_{rege}	S_{re}	S_{ge}	S_{rege}	
EmoLA (single task)	64.5	31.7	96.2	56.3	35.2	91.5	
EmoLA (multi-task)	64.3	32.0	96.3	54.7	33.9	88.6	

by a phase where only the LLM decoder is trained. LLaVA-1.5 enhances its performance by incorporating a two-layer MLP and adopting higher image resolutions. Shikra shares a similar architecture with LLaVA; however, it distinguishes itself by fine-tuning both the projector and the LLM decoder during its training phases. Similarly, MiniGPT4-V2, while architecturally similar to LLaVA, employs higher-resolution images to improve visual perception and aggregates every four neighboring visual tokens into a single token to optimize training efficiency. mPLUG-Owl2 introduced a visual abstractor module which aggregates the information from visual tokens by learnable queries. Furthermore, mPLUG-Owl2 leverarged a modality-adaptive module for facilitating multi-modal inputs to a shared semantic space for enabling modality collaboration.

D.3 Multi-task performance on FABA-Bench

We also perform experiments to evaluate the multi-task performance of EmoLA on our FABA-Bench. Specifically, we train our EmoLA using dataset from two tasks instead of the individual dataset. As shown in Tab. 13, the performance on Emotion of EmoLA under multi-task setting is almost the same to the single task setting. While for AU task, the performance of multi-task EmoLA drops com-

Table 14: The location of prior token.

Methods	Emotion					
	S_{re}	S_{ge}	S_{rege}	S_{re}	S_{ge}	S_{rege}
EmoLA (prior token before visual tokens)	63.5	32.1	95.6	55.4	34.3	89.7
EmoLA (prior token after visual tokens)	64.5	31.7	96.2	56.3	35.2	91.5

pared to the single task version. We assume it's because the emotion recognition task will somehow affect the AU recognition performance by making model focus more on the general emotion recognition. Moreover, the multi-task is harder than the single task, which may also decrease the performance on AU task.

D.4 Location of prior token

We also investigate the location of prior token in EmoLA in Tab. 14. From the results, we can observe that if we put the prior token before the visual tokens the recognition performance on two tasks will decrease. We assume this will affect the visual tokens' representation due to the causal mask in the decoder. As a result, we put the prior token after the visual tokens.

E More generation cases

E.1 Successful cases

We present some successful cases of EmoLA's prediction on our FABA-Instruct dataset in Fig. 11 and Fig. 12, respectively.

E.2 Failed cases

We also present some failed cases of EmoLA's prediction on our FABA-Instruct dataset in Fig. 13a and Fig. 13b, respectively.

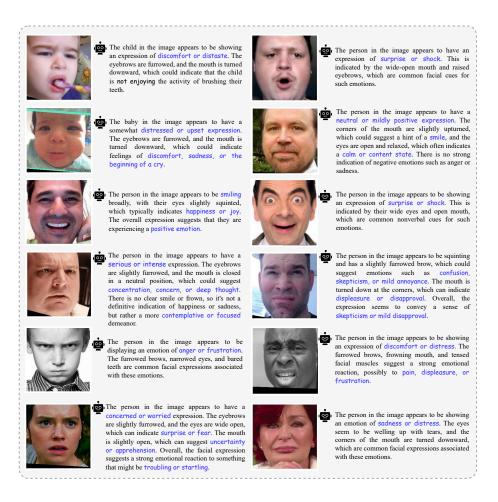
- "What is the emotion of this face?",
- "Tell me about the emotion on this face",
- "Label the emotion of this face?",
- "What emotions do you perceive in this face?",
- "How does the emotion on this face make you feel?",
- "Can you tell me the emotion displayed in this image?",
- "What emotional response does this face evoke?",
- "Can you determine the emotion portrayed in this image?",
- "Provide a detailed interpretation of the emotion on this face",
- "What can you say about the emotional expression in this image?",
- "Can you analyze the emotion depicted in this face?",
- · "What emotions come to mind when you see this face?",
- "How do you perceive the emotional tone of this face?",
- "Can you provide insights into the emotions conveyed in this image?",
- · "What emotions do you think this face is expressing?",
- "How does the emotion in this image resonate with you?",
- · "Can you articulate the emotional context of this face?",
- "What is your interpretation of the emotions displayed by this person?",
-

(a) Examples of emotion instructions.

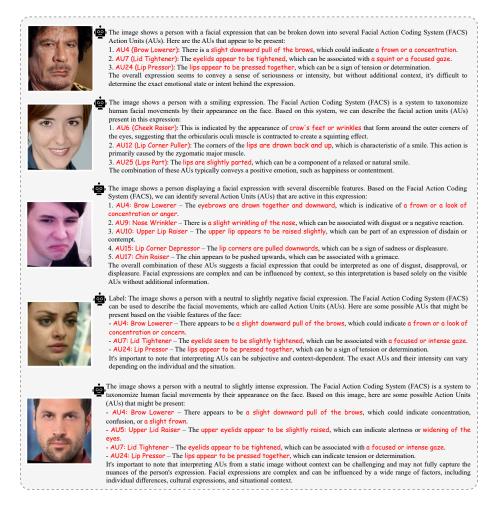
- "What is the action units of this face?"
- "List all the action units of this face."
- "What action units can be identified in this face?"
- "Please enumerate the action units observed in this particular face?"
- "Provide a breakdown of the action units present on this face?"
- "Identify the action units exhibited by this face."
- · "List the action units that characterize this face."
- "Can you specify the action units found in this face?"
- "Enumerate the action units displayed by this face."
- "Describe the action units visible on this face."
- "What action units are evident in this face?"
- "List the action units observable in this face."
- "Give a detailed list of action units for this face."
- "What are the action units present in this face?"
- "List all the action units of this face."
- "Please enumerate the action units observed in this particular face?"
- "What specific facial movements correspond to action units of this image?"
- "Provide a detailed breakdown of the action units visible in this face."
-

(b) Examples of AU instructions.

Fig. 10: The instructions used in FABA-Instruct.



 ${\bf Fig.\,11:}\ {\bf Examples\ of\ EmoLA's\ successful\ prediction\ on\ emotion\ task}.$



 $\textbf{Fig.\,12:} \ \ \textbf{Examples of EmoLA's successful prediction on AU task}.$



(b) Examples of EmoLA's failed prediction on AU task.

Fig. 13: Failed cases of EmoLA on FABA-Instruct.