

# ContextFace: Generating Facial Expressions from Emotional Contexts

Min-jung Kim  
 Dept. of Computer Science  
 and Engineering  
 Korea University  
 LG Electronics  
 South Korea

kmj0606@korea.ac.kr

Minsang Kim  
 Dept. of Computer Science  
 and Engineering  
 Korea University  
 SK Telecom  
 South Korea

kmswin1@korea.ac.kr

Seung Jun Baek\*  
 Dept. of Computer Science  
 and Engineering  
 Korea University  
 South Korea

sjbaek@korea.ac.kr

## Abstract

The task of generating 3D facial expressions given various situational contexts is important for applications such as virtual avatars or human-robot interactions. The task is, however, challenging not only because it requires a comprehensive understanding of emotion, expression and contexts, but also there rarely are datasets to support the task. We propose ContextFace, a Multi-modal Large Language Model (MLLM) fine-tuned to generate 3D facial expressions depending on complex situational contexts. To overcome the lack of datasets, we perform a context augmentation to existing emotion recognition datasets; we generate plausible situations and quotes from images and emotions to annotate the dataset. Next, we perform visual instruction tuning of MLLMs on context-augmented datasets to boost their capability of visual synthesis from emotions. Experiments show a superior performance of ContextFace in the zero-shot evaluation of contextual emotion recognition. A qualitative evaluation shows that our method generates expressions consistent with diverse contexts and performs complex emotion reasoning, e.g., speculative generation of expressions of occluded faces through interactive prompting. Code and data are available at [https://github.com/minjung98/ContextFace\\_.git](https://github.com/minjung98/ContextFace_.git).

## 1. Introduction

The interpretation of facial expressions inherently depends on the situational context. For instance, a smiling face can convey different emotions – from genuine joy to forced politeness – depending on the situational context. A combined understanding of expression, emotion, and context has recently become important, particularly in the field of human-

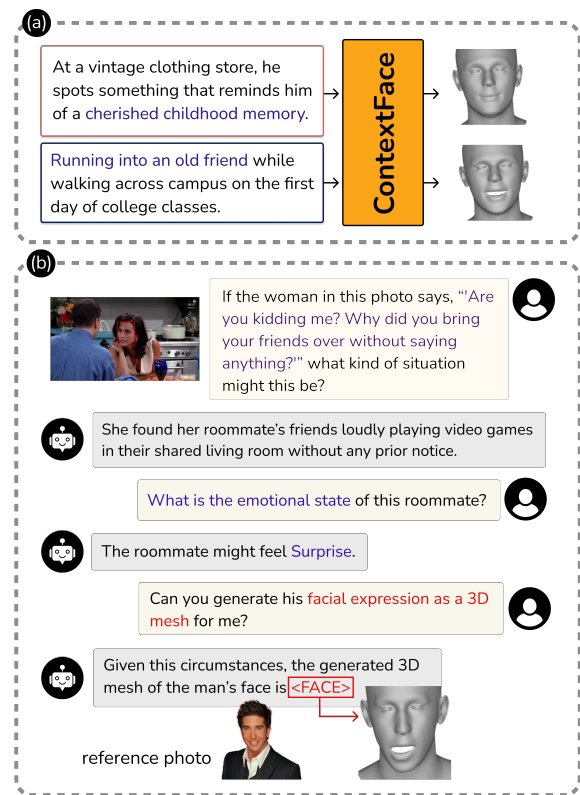


Figure 1. (a) **Context-to-face generation.** ContextFace generates subtly different expressions for the same emotion “Happy”. (b) **Speculative generation of expressions.** The emotion of woman with the image alone is ambiguous. But with her quote, ContextFace infers her emotion and generate plausible situations. Moreover, ContextFace speculatively generates the man’s expression whose face is not visible. An embarrassed face is generated from special token <FACE> through an interactive prompting. The output is FLAME parameters [19] which is later matched to the reference photo using INFERNO [7] for the final mesh.

\*Corresponding Author

robot interaction [20, 33, 39, 41]. Prior works on Facial Expression Recognition [32, 36, 40, 44] focused on facial feature analysis, but underexplored the richness of contextual information. Recently, Multimodal Large Language Models (MLLMs) enabled interpretable analysis of expressions by explaining how facial features indicate specific emotional states [21, 24, 46]. MLLMs can perform context-aware emotion recognition [4, 11, 17, 45, 47] to logically infer emotions by sophisticated reasoning, opening new avenues for emotion analysis from multi-modal cues.

While the aforementioned works focus on predicting emotions or textual explanations, we consider a task of *generating 3D facial expressions conditional on complex situational contexts*. The task has a wide range of applications, e.g., interactive visual assistants [25, 29] or virtual avatars [3, 13], where one needs to aptly synthesize 3D faces that well-reflect emotions arising from various situations. It is challenging to generate expressions capturing the subtlety of situational contexts. Consider two situations: “unexpectedly meeting an old friend on the street” and “cherishing childhood memories”. Both will be categorized as “happy” in traditional emotion recognition; however, the faces of the persons in the contexts will subtly differ, although both are expected to be “smiley faces”. The task requires reasoning capabilities to derive proper emotion and expressions from situational contexts. However, a lack of datasets that encompass expressions, emotions, and contexts makes it difficult to develop multi-modal models.

In this paper, we propose ContextFace, an MLLM that integrates contextual understanding and emotional reasoning to generate facial expressions. ContextFace is able to generate proper expressions matched to details in the contexts (Fig. 1(a)). Both situations imply emotion “happy”, but leading to the generation of different faces. Moreover, ContextFace can perform emotional reasoning in complex contexts. In Fig. 1(b), the model is prompted to 1) generate a plausible situation of the conversation; 2) infer the expression of the man whose face is not visible. Although the emotion of the woman with the image alone is ambiguous, ContextFace can infer her emotion from her quote, create reasonable situations and *speculatively* generate a proper facial expression of the man through an interactive prompting.

The design of ContextFace is based on *visual instruction tuning* [29, 30]. We first create datasets for instruction tuning using *context augmentation*. Specifically, we augment existing emotion-labeled image datasets [9, 16] with two new annotations: situation descriptions and subject quotes. For the annotation, we leverage the capabilities of strong LLMs in grounding emotions in commonly understood social contexts. The datasets will be publicly released for the research in emotional analysis. Next, we construct instruction datasets derived from the newly augmented datasets, fine-tune an MLLM on those instructions. The MLLM

uses special token <FACE> associated with its hidden state along with the face projection module to estimate FLAME parameters [19]. This enables 3D face reconstruction with expressions based on reference photos using INFERNO [7]. Experiments show the superior performance of ContextFace in contextualized emotion recognition compared to existing MLLMs as well as its capability of emotional reasoning and face generation given various contexts, such as speculative generation of occluded faces through interactive prompting.

Our contribution is summarized as follows. (1) We create and publicly release two emotion datasets augmented with rich contextual information. (2) We propose ContextFace, an MLLM fine-tuned with visual instruction tuning for synthesizing proper expressions from complex situational contexts. (3) Experiments show ContextFace has superior emotional reasoning and recognition capabilities in both quantitative and qualitative aspects.

## 2. Related Work

**MLLMs and Instruction Tuning.** Enabling large language models to process data modalities beyond text has been a foundational area of research [18, 29, 42, 43]. Flamingo [1] proposed general-purpose large vision language models trained from a large-scale web corpus. BLIP [18] proposed an efficient VLM training strategy with Q-Former to bridge the image and text modality using cross-attention. LLaVA [29, 31] proposed a 2-stage visual instruction tuning method, which first aligned vision encoders with LLMs and then fine-tuned the VLM using chat-oriented datasets. The visual instruction tuning was successfully applied to various tasks. LISA [15] leverages the capabilities of LLMs for segmentation tasks that require complex reasoning. GLaMM [38] generates natural conversational texts with integrated pixel-level object segmentation masks in visual interactions. ChatPose [12] enables text and image inputs to generate 3D human body poses while supporting complex reasoning about posture. Inspired by those works, we apply visual instruction tuning for generation of facial expressions given emotional contexts.

**3D Face Reconstruction.** FLAME [19] is a high-quality 3D face reconstruction model from a single image, which is a parameterized framework for 3D face representation with separate components for identity, pose, and expression. EMOCA [6] proposed an emotion consistency loss to generate emotional expressions of significantly higher fidelity. Recent works such as EmoTalk [37], EMOTE [8] and EmoFace [26] proposed emotion-aware 3D face reconstruction methods that incorporate emotional expressions into 3D face representations by analyzing speech. However, integrating contextually-aware emotional expressions into 3D face representations remains an underexplored task.

**Multimodal Emotion Understanding.** DialogueLLM [47] proposed a finetuning of LLMs with multimodal emotional

dialogues from texts and videos, and released datasets on emotion recognition conversations. There are works [10, 45] which studied reasoning emotions within the situational context, and showed that recognizing accurate visual markers such as bounding boxes could significantly boost the performance. Emotion LLaMA [4] has presented an MLLM that infers human emotions by integrating multiple cues including facial expressions analyzed through action units, utterances, audio tones, and visual context descriptions. Explainable Multimodal Emotion Recognition (EMER) [23] was proposed to improve the reliability and accuracy of emotion recognition by generating explanations for their predictions. While those approaches leveraged multimodality to recognize and reason about emotions, our research takes a step forward by introducing a novel task of generating facial expressions from emotional reasoning based on complex situational contexts.

### 3. Method

We propose ContextFace, a Multimodal Large Language Model (MLLM) that integrates contextual understanding, emotional reasoning and facial expression generation. We aim to achieve this with *visual instruction tuning* inspired by recent works [12, 15, 38]. However, there rarely exists datasets that integrate context, emotion and facial expressions together. Our strategy to overcome the data scarcity is to perform *contextual augmentation* to existing emotion datasets assisted by LLMs. The augmented datasets enable us to create an instruction-tuning pipeline for an integrated understanding of contexts and emotions to generate appropriate facial expressions. With newly defined instruction tasks and datasets, we propose an efficient MLLM architecture for 3D facial mesh generation. In the following sections, we describe our dataset augmentation approach using LLMs (Sec. 3.1), visual instruction tuning (Sec. 3.2), and the proposed architecture (Sec. 3.3).

#### 3.1. LLM-Assisted Context Augmentation to Emotion Dataset

**Datasets.** We construct emotion datasets augmented with contextual annotations. We publicly release two datasets: SFEW-C (Static Facial Expressions in the Wild with Context) and CAER-S-C (Context-Aware Emotion Recognition - Static with Context). Our datasets are based on and enhance SFEW [9] and CAER-S [16], where these original datasets contain samples of image (input) and emotion (label) pairs. We chose SFEW and CAER-S, because they comprise frame captures of human subjects from movies and TV shows, offering a good balance between contextual background information and clear facial expressions of the subjects. We augment each sample of SFEW and CAER-S with contextual information; specifically, situation and quote for the image-emotion pair, e.g., see Fig. 2(a). With

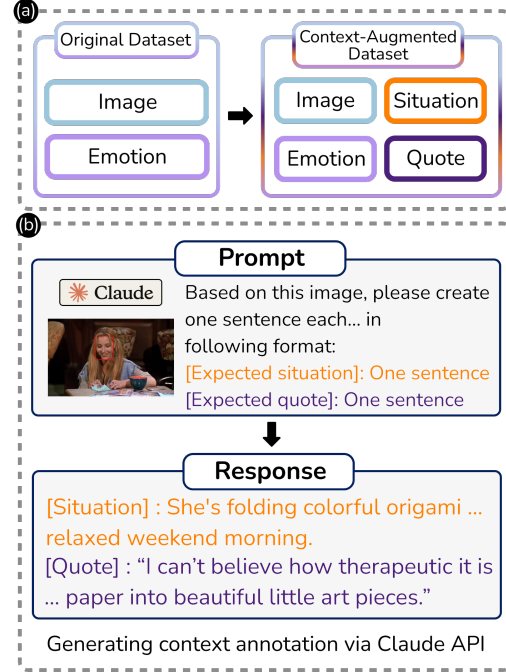


Figure 2. **Contextual Augmentation to Emotion Datasets.** (a) The original dataset is augmented with contexts, plausible situations and quotes, consistent with the image and emotion. (b) The context annotations are generated through prompting the Claude API with images and emotion labels.

the contextual annotations, our datasets will be useful for training MLLMs to understand emotion, context and facial expressions.

**Contextual Annotation.** We generate the contextual annotations assisted by strong LLMs. We leverage the vast visual and textual knowledge of LLMs to infer contexts from image and emotion. Given the original dataset’s images and emotion labels as input, we prompted LLMs to generate expected situation and plausible quotes (Fig. 2(b)). The annotations are generated through the Claude-3-5-sonnet-20241022 API [2]. From the original datasets, we excluded several samples that did not meet our criteria, e.g., images without human subjects, those containing adult/inappropriate content rejected by the Claude API, etc. After applying our preprocessing criteria, our final released annotation dataset contains 42,196 training and 20,938 test samples for CAER-S-C, and 890 training and 431 test samples for SFEW-C.

#### 3.2. Visual Instruction Tuning

To train ContextFace, we construct instruction tuning datasets derived from the augmented datasets introduced in Sec. 3.1. We define two new tasks: Situation Generation and Expression Generation and build the associated instruction datasets, along with a generic VQA dataset [29] for the

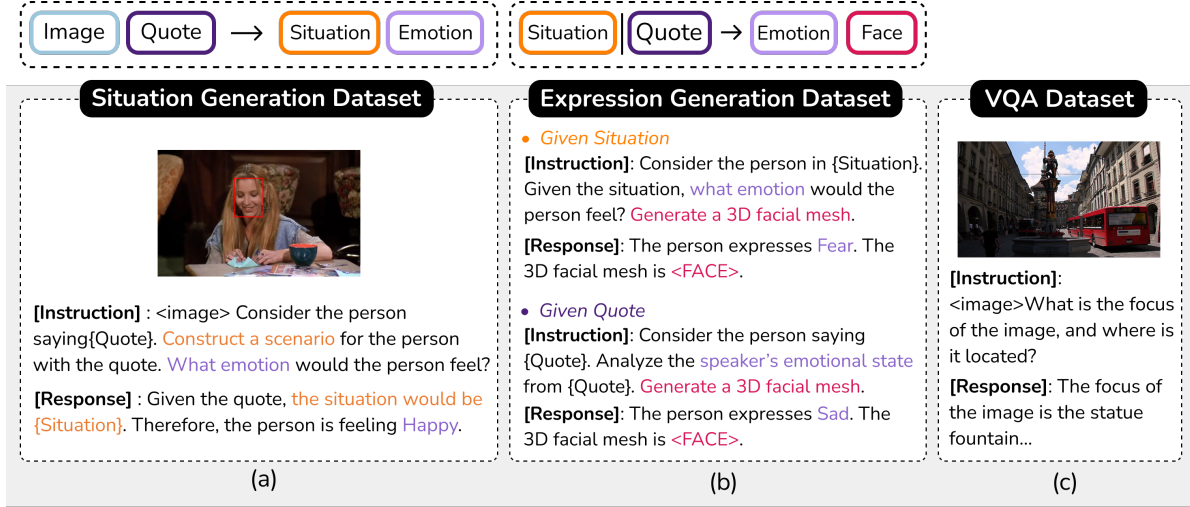


Figure 3. **Instruction datasets consisting of three types:** (a) Situation Generation Dataset, (b) Expression Generation Dataset, and (c) Visual Question Answering Dataset.

instruction tuning of ContextFace.

**Situation Generation Instructions.** Consider the augmented dataset in Sec. 3.1 consisting of image, emotion, situation and quote. From the dataset, we create a task called *Situation Generation* which takes image and quote as input, and predicts situation and emotion as response (see Fig. 3(a)). This task enables a model to learn associations between contextual cues and emotions.

We create a instruction tuning dataset for Situation Generation task as follows. We prepare training pairs  $\{x_{img}, x_{txt}\}$  and their corresponding text output  $y_{txt}$  from the augmented dataset. Here,  $x_{img}$  denotes the input image with a red bounding box marking the target person.  $x_{txt}$  is a text prompt which asks the model to create a scenario from the input quote and predict the emotion from predefined categories.  $y_{txt}$  denotes the target response describing a plausible situation that fits the quote, followed by the identified emotional state. We organized the data in an instruction-response as follows, with <image> denoting the position of image tokens:

**INSTRUCTION:** <image> Think about a person who says “I can’t believe how therapeutic it is to turn these simple pieces of paper into beautiful little art pieces.” - construct the scenario that would result in this statement, and what emotion would they be feeling? Choose from: Angry, Sad, Surprise, Neutral, Fear, Happy, Disgust.

**RESPONSE:** She’s folding colorful origami while enjoying a relaxing afternoon coffee break at home. Therefore, the emotion this person is feeling is Happy.

**Expression Generation Instructions.** We create a task

called *Expression Generation* which takes context (situation or quote) as input, and predicts 3D facial expression and emotion as response (see Fig. 3(b)). With this task, the model learns to produce 3D facial expressions that correspond to its emotional analysis of given situations or quotes.

The training data comprises input  $\{x_{txt}\}$  and target response  $\{y_{txt}, \beta\}$  obtained from the augmented dataset.  $x_{txt}$  denotes a text instruction that describes a situation or quote and asks the model to generate an appropriate 3D facial representation.  $y_{txt}$  is the target response including both emotional inference and the special token called <FACE>. <FACE> is a learnable token from which 3D facial mesh will be predicted.  $\beta$  denotes the latent representation of the target 3D facial mesh. Specifically,  $\beta$  is a 100-dimensional FLAME expression parameters [19] extracted from the face in the ground truth image using INFERNO [7]. We structure the data in the following instruction-response format:

**INSTRUCTION:** She’s folding colorful origami while enjoying a relaxing afternoon coffee break at home. From the described situation, infer the person’s emotional state and generate an appropriate 3D facial mesh that captures the feeling.

**RESPONSE:** Considering the situation, the person expresses Happy. The 3D facial mesh is <FACE>.

**INSTRUCTION:** “I can’t believe how therapeutic it is to turn these simple pieces of paper into beautiful little art pieces.” Please analyze the speaker’s emotional state from this quote and generate a corresponding 3D facial mesh.



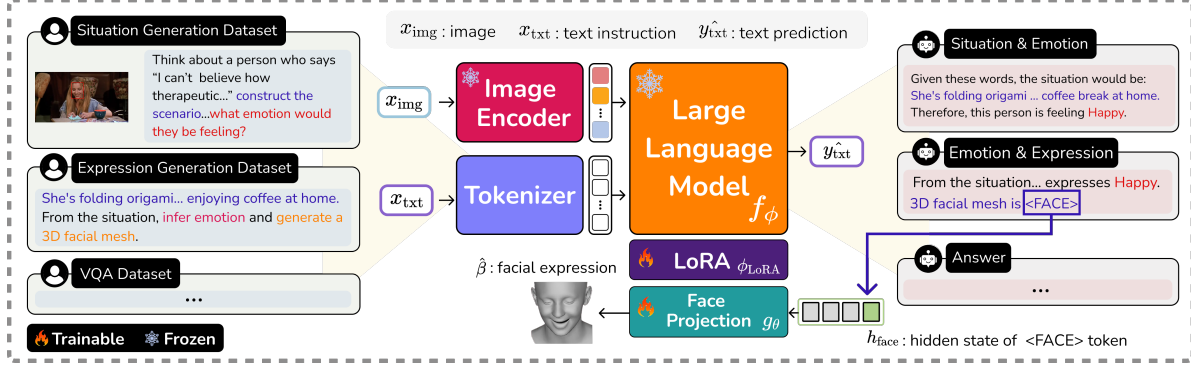


Figure 4. **Pipeline of ContextFace.** The instruction datasets described in Fig. 3 are used as input. Images and text are processed by the image encoder and tokenizer respectively, and are fed into the LLM to produce text output. For the expression generation task, the hidden state of the <FACE> token in the text output is processed through face projection layers to generate the coefficient vector of facial expressions. The coefficient vector is then used to generate the final facial meshes.

**RESPONSE:** Considering the quote, the person expresses **Happy**. The 3D facial mesh is **<FACE>**.

**Visual Question Answering Dataset.** To preserve the model’s original capabilities in Visual Question Answering during the training on task-specific data, we incorporated VQA dataset (Fig. 3(c)), specifically LLaVA-Instruct-150K [29] which is a visual-language instruction dataset generated by GPT-4, into our training pipeline.

### 3.3. Architecture

**Model Design.** As illustrated in Fig. 4, our model consists of two main components: a multi-modal large language model denoted as  $f_\phi$  and a face projection model  $g_\theta$ . Text instructions  $\mathbf{x}_{\text{txt}}$  serve as inputs to the multi-modal LLM  $f_\phi$ , resulting in text outputs  $\hat{\mathbf{y}}_{\text{txt}}$  given by

$$\hat{\mathbf{y}}_{\text{txt}} = f_\phi(\mathbf{x}_{\text{txt}}). \quad (1)$$

We employed LoRA fine-tuning [14] with trainable parameters  $\phi_{\text{LoRA}}$  to efficiently adapt the model to our specific task while preserving the generation capabilities of pre-trained vision-language models. If the input prompt requests the generation of 3D facial mesh, the multi-modal LLM  $f_\phi$  is trained to include <FACE> token in its text outputs  $\hat{\mathbf{y}}_{\text{txt}}$ . The hidden state of the multi-modal LLM  $f_\phi$  corresponding to <FACE> token, denoted by  $\mathbf{h}_{\text{face}}$ , is used to generate the 3D facial mesh.  $\mathbf{h}_{\text{face}}$  contains prompt-specific contextual information required to generate a facial expression.  $\mathbf{h}_{\text{face}}$  is subsequently passed through the face projection module  $g_\theta$  to obtain the expression coefficient  $\hat{\beta}$  given by

$$\hat{\beta} = g_\theta(\mathbf{h}_{\text{face}}). \quad (2)$$

We use a two-layer multi-layer perceptron (MLP) for  $g_\theta$ .

**Loss Functions.** The loss function is designed to optimize both text output  $\hat{\mathbf{y}}_{\text{txt}}$  and facial expression output  $\hat{\beta}$ . Firstly,

the loss  $\mathcal{L}_{\text{txt}}$  for the text generation is the cross entropy loss between the predicted response  $\hat{\mathbf{y}}_{\text{txt}}$  and the ground truth  $\mathbf{y}_{\text{txt}}$ :

$$\mathcal{L}_{\text{txt}} = \text{CE}(\mathbf{y}_{\text{txt}}, \hat{\mathbf{y}}_{\text{txt}}). \quad (3)$$

Secondly, for the facial expression loss denoted by  $\mathcal{L}_{\text{face}}$ , we adopted the L2 loss as the objective function for our facial expression component. The loss compares L2 distance between the predicted  $\hat{\beta}$  and ground truth  $\beta$  coefficients given by

$$\mathcal{L}_{\text{face}} = \|\beta - \hat{\beta}\|_2^2. \quad (4)$$

where  $\|\cdot\|_2$  denotes the L2-norm of a vector. The overall loss  $\mathcal{L}$  consists of text generation loss  $\mathcal{L}_{\text{txt}}$  and facial expression  $\mathcal{L}_{\text{face}}$ , which is given by

$$\mathcal{L} = \mathcal{L}_{\text{txt}} + \lambda \cdot \mathcal{L}_{\text{face}}. \quad (5)$$

where  $\lambda$  is a weight that balances the losses.

**Training Process.** ContextFace is jointly trained on three instruction datasets: Situation Generation, Expression Generation, and VQA, in an end-to-end manner. Specifically, the model is jointly trained with batches from different tasks, where each batch contained samples from a single dataset randomly selected according to fixed ratios. Both the token embedding matrix and LLM prediction head are set to be trainable to handle the newly added <FACE> token. Also, the face projection  $g_\theta$  is a trainable module to learn the mapping from  $\mathbf{h}_{\text{face}}$  to facial expression parameters  $\hat{\beta}$ . Due to the inherently smaller magnitude of the facial expression loss values, we applied a weighting factor  $\lambda$  of 10 to this term in the overall loss function during training.

## 4. Experiment

### 4.1. Implementation Details

**Datasets and Baselines.** We use two datasets SFEW [9] and CAER-S [16] which contain captures of movies

Method	Hap	Sad	Neu	Ang	Sur	Dis	Fea	UAR	WAR
<i>with situation</i>									
BLIP-13B [5]	72.41	55.45	50.92	81.82	21.21	64.86	42.62	51.11	51.09
llava-1.5-13B [29]	<b>96.50</b>	<u>93.96</u>	<u>90.45</u>	77.78	73.33	78.05	71.43	81.42	84.98
LLaVA-NEXT-13B [31]	93.62	93.51	88.31	<u>91.57</u>	69.90	<b>85.71</b>	<b>86.27</b>	<u>86.88</u>	<u>89.73</u>
Qwen2.5-VL-7B [43]	<u>95.10</u>	89.44	86.96	89.04	<u>73.79</u>	<u>79.17</u>	84.00	86.08	88.33
ContextFace (ours)	94.20	<b>95.89</b>	<b>92.77</b>	<b>94.94</b>	<b>85.45</b>	<b>85.71</b>	<u>84.31</u>	<b>90.13</b>	<b>90.68</b>
<i>with quote</i>									
BLIP-13B [5]	26.51	21.95	36.44	35.79	0.00	16.00	40.00	26.35	24.52
llava-1.5-13B [29]	<u>93.53</u>	<u>93.15</u>	83.22	86.36	<u>80.73</u>	74.42	78.00	84.12	85.66
LLaVA-NEXT-13B [31]	90.23	90.51	77.55	91.57	78.50	<u>83.72</u>	71.32	84.01	85.56
Qwen2.5-VL-7B [43]	90.65	<b>97.96</b>	<b>91.72</b>	<u>92.99</u>	80.70	76.36	<u>90.32</u>	<u>90.20</u>	<u>91.33</u>
ContextFace (ours)	<b>96.60</b>	92.21	<u>84.93</u>	<b>94.48</b>	<b>91.74</b>	<b>95.65</b>	<b>92.78</b>	<b>93.21</b>	<b>94.89</b>

Table 1. **Zero-shot emotion recognition performance in F1 scores on SFEW-C Dataset.** The input data for upper (resp. lower) table is image-situation (resp. image-quote) pair. Emotion categories: Hap (Happy), Sad (Sad), Neu (Neutral), Ang (Angry), Sur (Surprise), Dis (Disgust), Fea (Fear). UAR: Unweighted Average Recall, WAR: Weighted Average Recall.

and TV shows and are suitable for studying contextual emotions. For training and evaluation, we use context-augmented versions: SFEW-C and CAER-S-C as introduced in Sec. 3.1. ContextFace is trained with CAER-S-C. We compare ContextFace with state-of-the-art MLLM models: BLIP-13B [5], LLaVA-v1.5-13B [27], LLaVA-NEXT-13B [31] and Qwen2.5-VL-7B [43]. ContextFace and all the baseline models are evaluated in a zero-shot manner on SFEW-C. In addition, we provide an evaluation on emotion datasets derived from MER2023 [4, 22] datasets: due to space constraints, readers are referred to Sec. 4 of Supplementary Materials.

**Network Architecture.** We employed LLaVA-llama2-13B [28] with CLIP ViT-L/14 (336px) as our multimodal backbone. The face projection module is a two-layer MLP with input, hidden layer and output sizes given by 5120, 5120, and 100 respectively.

## 4.2. Quantitative Results

**Contextual Emotion Recognition.** Tab. 1 shows the performance of contextual emotion recognition. Specifically, the task is to classify images accompanied by one of two context types, situation or quote, into seven categories of emotions. ContextFace achieves superior performance in the majority of emotion categories for both situation and quote contexts, with the highest UAR (Unweighted Average Recall) and WAR (Weighted Average Recall) scores in each contextual setting (90.13%, 90.68% with situation context and 93.21%, 94.89% with quote context). Notably, ContextFace excels at detecting emotions like Disgust and Fear which typical models for emotion recognition struggle to identify. This superior performance suggests that ContextFace has developed a sophisticated understanding of the

Baseline	L2 ↓	FD ↓
<i>situation</i>		
Random	0.34± 0.0102	1.98 ± 0.2263
Mean	0.14± 0.0047	14.38 ± 0.4688
<b>Ours</b>	<b>0.07± 0.0055</b>	<b>0.81± 0.0711</b>
<i>quote</i>		
Random	0.31± 0.0093	2.42 ± 0.1982
Mean	0.12 ± 0.0042	11.76 ± 0.4178
<b>Ours</b>	<b>0.10± 0.0048</b>	<b>1.95± 0.1305</b>

Table 2. **Prediction Performance of Facial Expression.** We measure the errors in the predicted facial expression given two context types: quotes and situations. There are two baselines. **Random:** distance between the prediction and a randomly selected sample from the test set (excluding the ground truth). **Mean:** distance between the prediction and the average of samples in the test set. **Ours:** distance between the prediction and the ground truth.

complex relationship between context and emotion through the proposed framework of visual instruction tuning.

**Facial Expression Prediction.** We assess the performance of predicting facial expression coefficient  $\beta$  for given situation or context. The expression-context pairs in SFEW-C are used for evaluation. We measure the error in the prediction in terms of the distance between 100-dimensional coefficient vectors of facial expressions in the L2 and FD (Fréchet Distance) metrics. L2 measures the error in the L2 distance between the coefficient vectors. FD measures the statistical distance between the distributions of coefficient vectors. Since there exist no other methods for contextual expression generation, we perform a self-evaluation

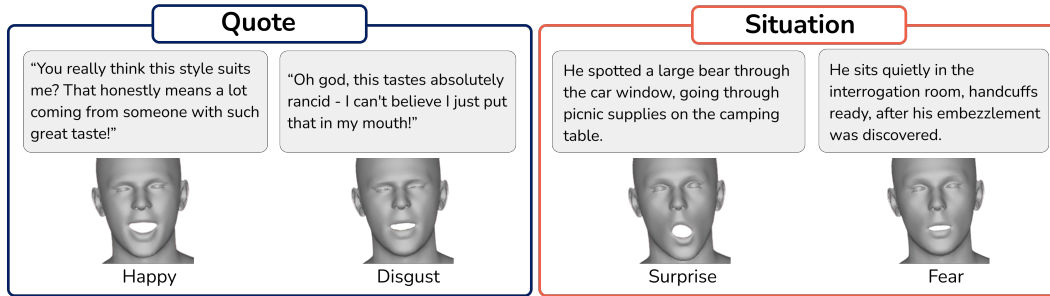


Figure 5. **Contextual Alignment of Emotional Expressions.** Generation results of facial expressions in response to quotes and situations.

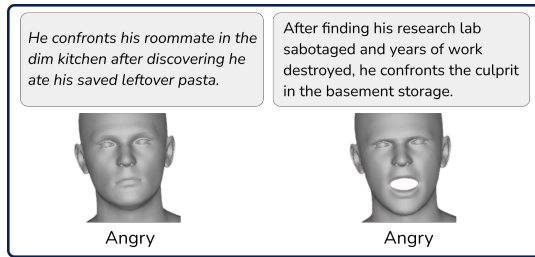


Figure 6. **Contextual Variation within Emotions.** Contextual variations of angry facial expressions in different situations.

adopting the multiple axes-based evaluation framework discussed in previous studies [34, 35] as follows. Tab. 2 shows the comparisons with Random and Mean baselines where we compare the distance between the prediction and a randomly selected sample from test set (Random), or the centroid of the test set (Mean), or the ground truth (Ours). Our model demonstrates a significant improvement over the Random baseline, suggesting that it has successfully learned the relationship between context and facial expressions. Furthermore, our model demonstrates substantially superior FD scores compared to the Mean baseline, preserving the distinctive elements of specific facial expressions to effectively capture fine facial feature details.

### 4.3. Qualitative Results

**Context-to-Emotional Expression Mapping.** As shown in Fig. 5, ContextFace can generate facial expressions aligned with each emotion by taking textual contexts, such as quotes or situations, as input.

**Intra-Emotion Expressional Diversity.** In Fig. 6, we observe that ContextFace is capable of generating different expressions within the same emotion depending on the situational context. Facial expressions involving anger were generated in different situational contexts. The left image displays a milder expression in response to a minor annoyance, while the right image shows a significantly more intense expression of rage in response to an infuriating situation.

**Speculative Expression Generation through Emotional**

**Reasoning.** We show that our model can generate logically valid situations when given a quote for ambiguous facial expressions. Additionally, based on these generated situations, the model can create appropriate facial expressions. In Fig. 7, we observe that our model correctly infers the emotions of the speaker, even with the same image and different quotes. Moreover, the model can infer the emotional state or facial expression of not only the speaker, but also the other person. As shown in the example, while the quote comes from a woman, the model can *speculatively* generate the facial expression of her boyfriend through an interactive prompting. This suggests that our model has emotional reasoning capabilities to understand interpersonal relationships and infer corresponding emotions or facial expressions.

**Facial expression generation for occluded faces.** As shown in Fig. 8, our model can infer the situation based on contexts (image and quote), and can generate plausible facial expressions, even when the face is intentionally and completely occluded. Our model successfully infers situations and generates appropriate expressions, demonstrating it has comprehensively learned the relationships among context, emotions, and facial expressions.

### 4.4. Ablation study

**Impact of Dataset Combinations.** As shown in Fig. 3, we used three datasets to jointly train ContextFace, i.e., two context-augmented datasets: Situation Generation Dataset, Expression Generation Dataset and one VQA Dataset. We measured the model’s ability to infer emotions with or without context-augmented datasets where VQA dataset is used by default. The results show substantial performance gains when both context datasets are incorporated. This demonstrates the importance of training with context-augmented data for emotion recognition.

**Impact of Model Scale.** The LLaVA-llama2-13B model demonstrated higher performance in both UAR and WAR compared to LLaVA-llama2-7B, as shown in Tab. 4. This suggests that as model size increases, emotional reasoning and expression generation capabilities improve. However,

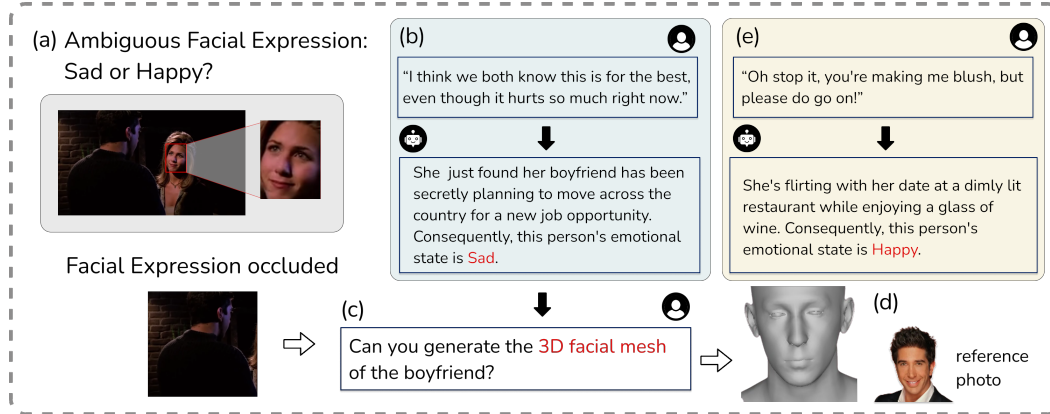


Figure 7. **Speculative Expression Generation through Emotional Reasoning.** (a) The expression in the input image is ambiguous. (b) Given the woman’s quote, our model can generate a plausible situation from the image along with the proper emotion. (c) The model is prompted to **speculate** the expression of the man whose face is not visible. (d) The model uses emotional reasoning to generate a sad face which is matched to a reference photo using INFERNO [7]. (e) When prompted with the same input image and a different quote, the model correctly predicts the woman’s emotion as “Happy”.

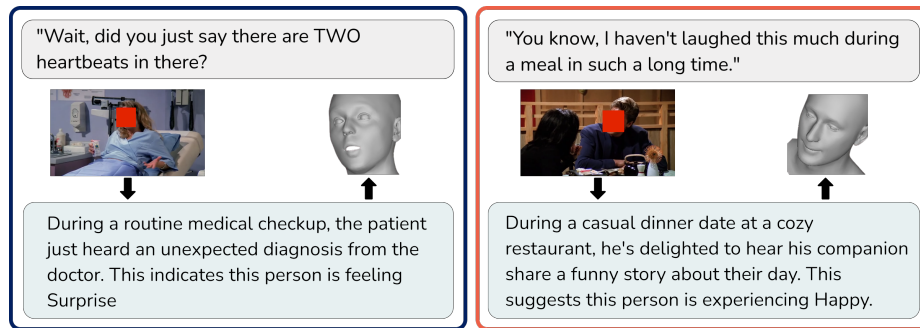


Figure 8. Generating situations and expressions from intentionally occluded faces and quotes.

Situation Generation	Expression Generation	UAR	WAR
×	×	68.26	70.53
×	✓	79.34	82.48
✓	×	88.33	89.40
✓	✓	<b>93.21</b>	<b>94.89</b>

Table 3. **Ablation study.** Performance comparison of emotion inference with different dataset combinations. ✓ (resp. ×) means the model is trained with (resp. without) the dataset.

the 7B model maintains competitive performance across both metrics, offering a computationally efficient alternative.

## 5. Conclusion

In this paper, we proposed ContextFace, a model that understands the deep association between contextual emotions and facial expressions and generates context-aligned facial expressions. Furthermore, we introduce new CAER-S-C

Method	UAR (%)	WAR
LLaVA-llama2-7B	89.71	90.56
LLaVA-llama2-13B	<b>93.21</b>	<b>94.89</b>

Table 4. **Ablation study.** Comparing emotion inference performance between LLaVA-llama2-13B and LLaVA-llama2-7B.

and SFEW-C datasets that enable the proposed contextual learning process. To our knowledge, we are the first to extend MLLMs beyond textual emotion analysis to direct generation of facial expressions, bridging the gap between language-driven emotion interpretation and visual synthesis. The proposed integration will become increasingly valuable as AI services evolve towards human-like interactions through emotionally aware and visually expressive communication.



## Acknowledgments

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korean government (MSIT) (2022R1A5A1027646), and the Institute of Information & Communications Technology Planning & Evaluation (IITP)-ICT Creative Consilience Program grant funded by the Korea government (MSIT) (IITP-2025-RS-2020-II201819).

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. 2
- [2] Anthropic. Claude-3.5-sonnet-20241022. <https://www.anthropic.com>, 2024. Large Language Model API. 3
- [3] Zehranaz Canfes, M Furkan Atasoy, Alara Dirik, and Pinar Yanardag. Text and image guided 3d avatar generation and manipulation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4421–4431, 2023. 2
- [4] Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Jingdong Sun, Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and Alexander Hauptmann. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning, 2024. 2, 3, 6
- [5] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 6
- [6] Radek Danecek, Michael J. Black, and Timo Bolkart. EMOCA: Emotion driven monocular face capture and animation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20311–20322, 2022. 2
- [7] Radek Daněček, Timo Bolkart, and Wojciech Zielonka. Inferno. <https://github.com/radekd91/inferno>, 2020. 1, 2, 4, 8
- [8] Radek Daněček, Kiran Chhatre, Shashank Tripathi, Yandong Wen, Michael Black, and Timo Bolkart. Emotional speech-driven animation with content-emotion disentanglement. In *SIGGRAPH Asia 2023 Conference Papers*, page 1–13. ACM, 2023. 2
- [9] Abhinav Dhall, Roland Goecke, Simon Lucey, and Tom Gedeon. Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark. In *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, pages 2106–2112, 2011. 2, 3, 5
- [10] Yasaman Etesam, Özge Nilay Yalçın, Chuxuan Zhang, and Angelica Lim. Contextual emotion recognition using large vision language models. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4769–4776. IEEE, 2024. 3
- [11] Yasaman Etesam, Özge Nilay Yalçın, Chuxuan Zhang, and Angelica Lim. Contextual emotion recognition using large vision language models, 2024. 2
- [12] Yao Feng, Jing Lin, Sai Kumar Dwivedi, Yu Sun, Priyanka Patel, and Michael J. Black. ChatPose: Chatting about 3d human pose. In *CVPR*, 2024. 2, 3
- [13] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Vid2avatar: 3d avatar reconstruction from videos in the wild via self-supervised scene decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12858–12868, 2023. 2
- [14] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 5
- [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023. 2, 3
- [16] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and Kwanghoonn Sohn. Context-aware emotion recognition networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2019. 2, 3, 5
- [17] Yuxuan Lei, Dingkan Yang, Zhaoyu Chen, Jiawei Chen, Peng Zhai, and Lihua Zhang. Large vision-language models as emotion recognizers in context awareness, 2024. 2
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. 2
- [19] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans, 2017. 1, 2, 4
- [20] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang, Ya Jing, Weinan Zhang, Huaping Liu, et al. Vision-language foundation models as effective robot imitators. In *The Twelfth International Conference on Learning Representations*. 2
- [21] Yifan Li, Anh Dao, Wentao Bao, Zhen Tan, Tianlong Chen, Huan Liu, and Yu Kong. Facial affective behavior analysis with instruction tuning, 2024. 2
- [22] Zheng Lian, Haiyang Sun, Licai Sun, Kang Chen, Mingyu Xu, Kexin Wang, Ke Xu, Yu He, Ying Li, Jinming Zhao, Ye Liu, Bin Liu, Jiangyan Yi, Meng Wang, Erik Cambria, Guoying Zhao, Björn W. Schuller, and Jianhua Tao. Mer 2023: Multi-label learning, modality robustness, and semi-supervised learning, 2023. 6
- [23] Zheng Lian, Licai Sun, Mingyu Xu, Haiyang Sun, Ke Xu, Zhuofan Wen, Shun Chen, Bin Liu, and Jianhua Tao. Explainable multimodal emotion reasoning. *CoRR*, 2023. 3
- [24] Zheng Lian, Haiyang Sun, Licai Sun, Jiangyan Yi, Bin Liu, and Jianhua Tao. Affectgpt: Dataset and framework for explainable multimodal emotion recognition, 2024. 2
- [25] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pre-training for vi-

- sual language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26689–26699, 2024. 2
- [26] Yihong Lin, Liang Peng, Jianqiao Hu, Xiandong Li, Wenxiong Kang, Songju Lei, Xianjia Wu, and Huang Xu. Emo-face: Emotion-content disentangled speech-driven 3d talking face with mesh attention, 2024. 2
- [27] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 6
- [28] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. LLaVA-Llama-2-13B-chat: Lightning Preview. <https://huggingface.co/liuhaotian/llava-llama-2-13b-chat-lightning-preview>, 2023. Accessed: 2025-07-31. 6
- [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 2, 3, 5, 6
- [30] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 2
- [31] 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. 2, 6
- [32] Shu Liu, Yan Xu, Tongming Wan, and Xiaoyan Kui. A dual-branch adaptive distribution fusion framework for real-world facial expression recognition. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 2
- [33] Zhentao Liu, Min Wu, Weihua Cao, Luefeng Chen, Jianping Xu, Ri Zhang, Mengtian Zhou, and Junwei Mao. A facial expression emotion recognition based human-robot interaction system. *IEEE CAA J. Autom. Sinica*, 4(4):668–676, 2017. 2
- [34] Evonne Ng, Hanbyul Joo, Liwen Hu, Hao Li, Trevor Darrell, Angjoo Kanazawa, and Shiry Ginosar. Learning to listen: Modeling non-deterministic dyadic facial motion, 2022. 7
- [35] Evonne Ng, Sanjay Subramanian, Dan Klein, Angjoo Kanazawa, Trevor Darrell, and Shiry Ginosar. Can language models learn to listen?, 2023. 7
- [36] Mang Ning, Albert Ali Salah, and Itir Onal Ertugrul. Representation learning and identity adversarial training for facial behavior understanding, 2024. 2
- [37] Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu Zhu, Hongyan Liu, Jun He, and Zhaoxin Fan. Emotalk: Speech-driven emotional disentanglement for 3d face animation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2023. 2
- [38] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M. Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S. Khan. Glamm: Pixel grounding large multimodal model. *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 3
- [39] Niyati Rawal and Ruth Maria Stock-Homburg. Facial emotion expressions in human-robot interaction: A survey. *International Journal of Social Robotics*, 14(7):1583–1604, 2022. 2
- [40] Arnab Kumar Roy, Hemant Kumar Kathania, Adhitiya Sharma, Abhishek Dey, and Md. Sarfaraj Alam Ansari. Resemotenet: Bridging accuracy and loss reduction in facial emotion recognition. *IEEE Signal Processing Letters*, pages 1–5, 2024. 2
- [41] Matteo Spezialetti, Giuseppe Placidi, and Silvia Rossi. Emotion recognition for human-robot interaction: Recent advances and future perspectives. *Frontiers in Robotics and AI*, 7:532279, 2020. 2
- [42] Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv preprint arXiv:2412.03555*, 2024. 2
- [43] Qwen Team. Qwen2.5-vl, 2025. 2, 6
- [44] Chengpeng Wang, Li Chen, Lili Wang, Zhaofan Li, and Xuebin Lv. Qcs: Feature refining from quadruplet cross similarity for facial expression recognition, 2025. 2
- [45] Alexandros Xenos, Niki Maria Foteinopoulou, Ioanna Ntinou, Ioannis Patras, and Georgios Tzimiropoulos. Vllms provide better context for emotion understanding through common sense reasoning, 2024. 2, 3
- [46] Bohao Xing, Zitong Yu, Xin Liu, Kaishen Yuan, Qilang Ye, Weicheng Xie, Huanjing Yue, Jingyu Yang, and Heikki Kälviäinen. Emo-llama: Enhancing facial emotion understanding with instruction tuning, 2024. 2
- [47] Yazhou Zhang, Mengyao Wang, Youxi Wu, Prayag Tiwari, Qiuchi Li, Benyou Wang, and Jing Qin. DialogueLLM: Context and emotion knowledge-tuned large language models for emotion recognition in conversations, 2024. 2