

ContextFace: Generating Facial Expressions from Emotional Contexts

Anonymous ICCV submission

Paper ID ****

Abstract

001 *The task of generating 3D facial expressions given various*
 002 *situational contexts is important for applications such as*
 003 *virtual avatars or human-robot interactions. The task is,*
 004 *however, challenging not only because it requires a com-*
 005 *prehensive understanding of emotion, expression and con-*
 006 *texts, but also there rarely are datasets to support the task.*
 007 *We propose ContextFace, a Multi-modal Large Language*
 008 *Model (MLLM) fine-tuned to generate 3D facial expressions*
 009 *depending on complex situational contexts. To overcome*
 010 *the lack of datasets, we perform a context augmentation to*
 011 *existing emotion recognition datasets; we generate plausi-*
 012 *ble situations and quotes from images and emotions to an-*
 013 *notate the dataset. Next, we perform visual instruction tun-*
 014 *ing of MLLMs on context-augmented datasets to boost its*
 015 *capability of visual synthesis from emotions. Experiments*
 016 *show a superior performance of ContextFace in the zero-*
 017 *shot evaluation of contextual emotion recognition. A qual-*
 018 *itative evaluation shows that our method generates expres-*
 019 *sions consistent with diverse contexts and performs complex*
 020 *emotion reasoning, e.g., speculative generation of expres-*
 021 *sions of occluded faces through interactive prompting.*

1. Introduction

023 The interpretation of facial expressions inherently depends
 024 on the situational context. For instance, a smiling face can
 025 convey different emotions – from genuine joy to forced po-
 026 liteness – depending on the situational context. A combined
 027 understanding of expression, emotion, and context has re-
 028 cently become important, particularly in the field of human-
 029 robot interaction [20, 31, 38, 40]. Prior works on Facial Ex-
 030 pression Recognition [30, 34, 39, 43] focused on facial fea-
 031 ture analysis, but underexplored the richness of contextual
 032 information. Recently, Multimodal Large Language Mod-
 033 els (MLLMs) enabled interpretable analysis of expressions
 034 by explaining how facial features indicate specific emo-
 035 tional states [21, 24, 45]. MLLMs can perform context-
 036 aware emotion recognition [4, 11, 17, 44, 46] to logically
 037 infer emotions by sophisticated reasoning, opening new av-

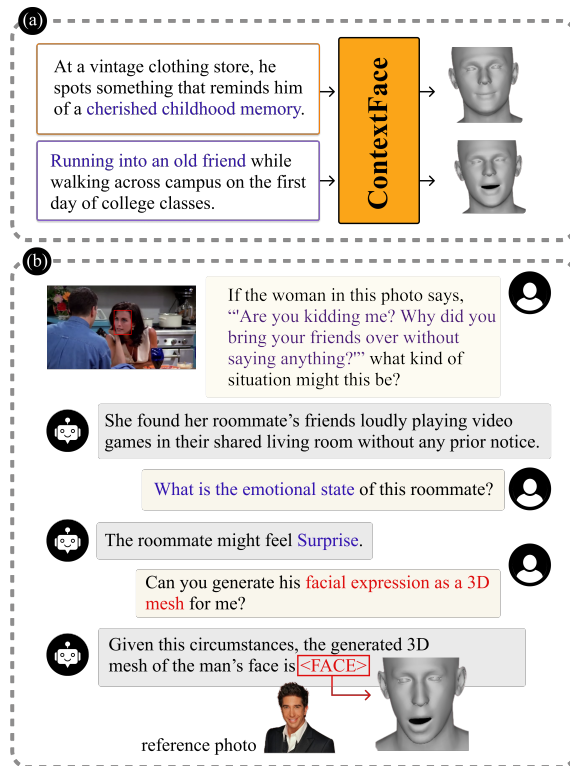


Figure 1. (a) **Context-to-face generation.** ContextFace generates subtly different expressions for the same emotion “Happy”. (b) **Speculative generation of expressions.** The emotion of woman with the image alone is ambiguous. But with her quote, ContextFace infers her emotion and generate plausible 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.

enues 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, 27] or virtual

038
039
040
041
042
043

044 avatars [3, 13], where one needs to aptly synthesize 3D
 045 faces that well-reflect emotions arising from various situ-
 046 ations. It is challenging to generate expressions capturing
 047 the subtlety of situational contexts. Consider two situa-
 048 tions: “unexpectedly meeting an old friend on the street”
 049 and “cherishing childhood memories”. Both will be catego-
 050 rized as “happy” in traditional emotion recognition; how-
 051 ever, the faces of the persons in the contexts will subtly dif-
 052 fer, although both are expected to be “smiley faces”. The
 053 task requires reasoning capabilities to derive proper emo-
 054 tion and expressions from situational contexts. However, a
 055 lack of datasets that encompass expressions, emotions, and
 056 contexts makes it difficult to develop multi-modal models.

057 In this paper, we propose ContextFace, an MLLM that
 058 integrates contextual understanding and emotional reason-
 059 ing to generate facial expressions. ContextFace is able to
 060 generate proper expressions matched to details in the con-
 061 texts (Fig. 1(a)). Both situations imply emotion “happy”,
 062 but leading to the generation of different faces. Moreover,
 063 ContextFace can perform emotional reasoning in complex
 064 contexts. In Fig. 1(b), the model is prompted to 1) gener-
 065 ate a plausible situation of the conversation; 2) infer the ex-
 066 pression of the man whose face is not visible. Although the
 067 emotion of the woman with the image alone is ambiguous,
 068 ContextFace can infer her emotion from her quote, create
 069 reasonable situations and *speculatively* generate a proper fa-
 070 cial expression of the man through an interactive prompting.

071 The design of ContextFace is based on *visual instruction*
 072 *tuning* [27, 28]. We first create datasets for instruction tun-
 073 ing using *context augmentation*. Specifically, we augment
 074 existing emotion-labeled image datasets [9, 16] with two
 075 new annotations: situation descriptions and subject quotes.
 076 For the annotation, we leverage the capabilities of strong
 077 LLMs in grounding emotions in commonly understood so-
 078 cial contexts. The datasets will be publicly released for the
 079 research in emotional analysis. Next, we construct instruc-
 080 tion datasets derived from the newly augmented datasets,
 081 fine-tune an MLLM on those instructions. The MLLM
 082 uses special token <FACE> associated with its hidden state
 083 along with the face projection module to estimate FLAME
 084 parameters [19]. This enables 3D face reconstruction with
 085 expressions based on reference photos using INFERNO [7].
 086 Experiments show the superior performance of ContextFace
 087 in contextualized emotion recognition compared to existing
 088 MLLMs as well as its capability of emotional reasoning and
 089 face generation given various contexts, such as speculative
 090 generation of occluded faces through interactive prompting.

091 Our contribution is summarized as follows. (1) We cre-
 092 ate and publicly release two emotion datasets augmented
 093 with rich contextual information. (2) We propose Con-
 094 textFace, an MLLM fine-tuned with visual instruction tun-
 095 ing for synthesizing proper expressions from complex situ-
 096 ational 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 lan-
 guage models to process data modalities beyond text has
 been a foundational area of research [18, 27, 41, 42].
 Flamingo [1] proposed general-purpose large vision lan-
 guage 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 us-
 ing cross-attention. LLaVA [27, 29] proposed a 2-stage vi-
 sual 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 [37] generates natural conver-
 sational texts with integrated pixel-level object segmenta-
 tion 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 genera-
 tion 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 expres-
 sion. EMOCA [6] proposed an emotion consistency loss
 to generate emotional expressions of significantly higher fi-
 delity. Recent works such as EmoTalk [35], EMOTE [8]
 and EmoFace [26] proposed emotion-aware 3D face recon-
 struction methods that incorporate emotional expressions
 into 3D face representations by analyzing speech. How-
 ever, it is an underexplored task to integrate the emotional
 expressions and complex situational contexts with ambigu-
 ous face images into 3D face representations.

Multimodal Emotion Understanding. DialogueLLM [46]
 proposed a finetuning of LLMs with multimodal emotional
 dialogues from texts and videos, and released datasets on
 emotion recognition conversations. There are works [10,
 44] which studied reasoning emotions within the situa-
 tional context, and showed that recognizing accurate vi-
 sual markers such as bounding boxes could significantly
 boost the performance. Emotion LLaMA [4] has presented
 an MLLM that infers human emotions by integrating mul-
 tiple cues including facial expressions analyzed through
 action units, utterances, audio tones, and visual context
 descriptions. Explainable Multimodal Emotion Recogni-
 tion (MER) [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

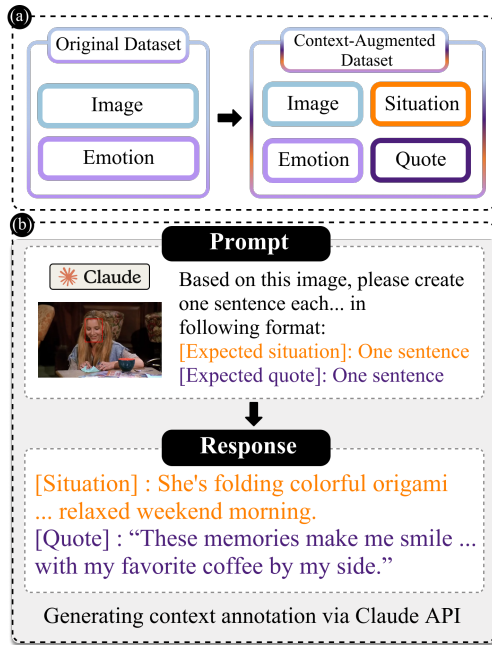


Figure 2. **Contextual Augmentation to Emotion Datasets.** (a) The original dataset is augmented with contexts, plausible situations and quotes, consistent with the image and emotion. (b) The context annotations are generated through prompting the Claude API with images and emotion labels.

149 research takes a step forward by introducing a novel task
150 of generating facial expressions from emotional reasoning
151 based on complex situational contexts.

152 3. Method

153 We propose ContextFace, a Multimodal Large Language
154 Model (MLLM) that integrates contextual understanding,
155 emotional reasoning and facial expression generation. We
156 aim to achieve this with *visual instruction tuning* inspired
157 by recent works [12, 15, 37]. However, there rarely exists
158 datasets that integrate context, emotion and facial expres-
159 sions together. Our strategy to overcome the data scarcity
160 is to perform *contextual augmentation* to existing emotion
161 datasets assisted by LLMs. The augmented datasets enable
162 us to create an instruction-tuning pipeline for an integrated
163 understanding of contexts and emotions to generate appro-
164 priate facial expressions. With newly defined instruction
165 tasks and datasets, we propose an efficient MLLM architec-
166 ture for 3D facial mesh generation. In the following sec-
167 tions, we describe our dataset augmentation approach using
168 LLMs (Sec. 3.1), visual instruction tuning (Sec. 3.2), and
169 the proposed architecture (Sec. 3.3).

3.1. LLM-Assisted Context Augmentation to Emotion Dataset 170 171

172 **Datasets.** We construct emotion datasets augmented with
173 contextual annotations. We publicly release two datasets:
174 SFEW-C (Static Facial Expressions in the Wild with Con-
175 text) and CAER-S-C (Context-Aware Emotion Recognition
176 - Static with Context). Our datasets are based on and en-
177 hance SFEW [9] and CAER-S [16], where these original
178 datasets contain samples of image (input) and emotion (la-
179 bel) pairs. We chose SFEW and CAER-S, because they
180 comprise frame captures of human subjects from movies
181 and TV shows, offering a good balance between contextual
182 background information and clear facial expressions of the
183 subjects. We augment each sample of SFEW and CAER-
184 S with contextual information; specifically, situation and
185 quote for the image-emotion pair, e.g., see Fig. 2(a). With
186 the contextual annotations, our datasets will be useful for
187 training MLLMs to understand emotion, context and facial
188 expressions.

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

3.2. Visual Instruction Tuning 204

205 To train ContextFace, we construct instruction tuning
206 datasets derived from the augmented datasets introduced in
207 Sec. 3.1. We define two new tasks: Situation Generation
208 and Expression Generation and build the associated instruc-
209 tion datasets, along with a generic VQA dataset [27] for the
210 instruction tuning of ContextFace.

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

218 We create a instruction tuning dataset for Situation
219 Generation task as follows. We prepare training pairs
220 $\{x_{img}, x_{txt}\}$ and their corresponding text output y_{txt}
221 from the augmented dataset. Here, x_{img} denotes the input image

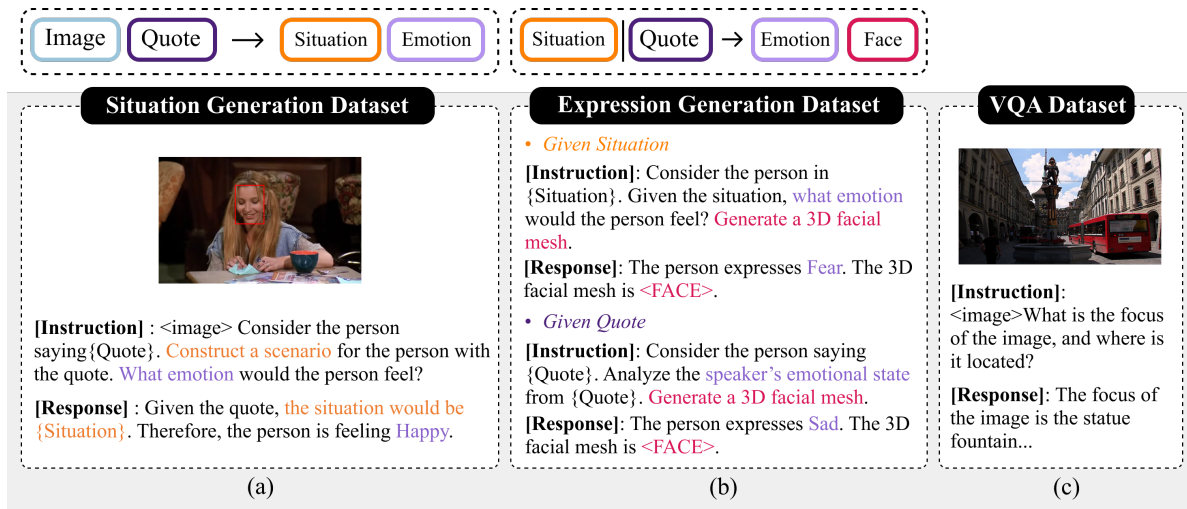


Figure 3. Instruction datasets consisting of three types: (a) Situation Generation Dataset, (b) Expression Generation Dataset, and (c) Visual Question Answering Dataset.

222 with a red bounding box marking the target person. x_{txt} is
 223 a text prompt which asks the model to create a scenario from
 224 the input quote and predict the emotion from predefined cat-
 225 egories. y_{txt} denotes the target response describing a plau-
 226 sible situation that fits the quote, followed by the identified
 227 emotional state. We organized the data in an instruction-
 228 response as follows, with <image> denoting the position
 229 of image tokens:

230 **INSTRUCTION:** <image> Think about a per-
 231 son who says “I can’t believe how therapeutic it
 232 is to turn these simple pieces of paper into beau-
 233 tiful little art pieces.” - construct the scenario that
 234 would result in this statement, and what emotion
 235 would they be feeling? Choose from: Angry, Sad,
 236 Surprise, Neutral, Fear, Happy, Disgust.
 237 **RESPONSE:** She’s folding colorful origami
 238 while enjoying a relaxing afternoon coffee break
 239 at home. Therefore, the emotion this person is feel-
 240 ing is Happy.

241 **Expression Generation Instructions.** We create a task
 242 called *Expression Generation* which takes context (situa-
 243 tion or quote) as input, and predicts 3D facial expression
 244 and emotion as response (see Fig. 3(b)). With this task, the
 245 model learns to produce 3D facial expressions that corre-
 246 spond to its emotional analysis of given situations or quotes.

247 The training data comprises input $\{x_{\text{txt}}\}$ and target re-
 248 sponse $\{y_{\text{txt}}, \beta\}$ obtained from the augmented dataset. x_{txt}
 249 denotes a text instruction that describes a situation or quote
 250 and asks the model to generate an appropriate 3D facial
 251 representation. y_{txt} is the target response including both
 252 emotional inference and the special token called <FACE>.
 253 <FACE> is a learnable token from which 3D facial mesh

will be predicted. β denotes the latent representation of the
 target 3D facial mesh. Specifically, β is a 100-dimensional
 FLAME expression parameters [19] extracted from the face
 in the ground truth image using INFERNO [7]. We struc-
 ture the data in the following instruction-response format:

254 **INSTRUCTION:** She’s folding colorful origami
 255 while enjoying a relaxing afternoon coffee break
 256 at home. From the described situation, infer the
 257 person’s emotional state and generate an appro-
 258 priate 3D facial mesh that captures the feeling.
 259 **RESPONSE:** Considering the situation, the per-
 260 son expresses Happy. The 3D facial mesh is
 261 <FACE>.

262 **INSTRUCTION:** “I can’t believe how therapeu-
 263 tic it is to turn these simple pieces of paper into
 264 beautiful little art pieces.” Please analyze the
 265 speaker’s emotional state from this quote and gen-
 266 erate a corresponding 3D facial mesh.
 267 **RESPONSE:** Considering the situation, the per-
 268 son expresses Happy. The 3D facial mesh is
 269 <FACE>.

270 **Visual Question Answering Dataset.** To preserve the
 271 model’s original capabilities in Visual Question Answering
 272 during the training on task-specific data, we incorporated
 273 VQA dataset (Fig. 3(c)), specifically LLaVA-Instruct-150K
 274 [27] which is a visual-language instruction dataset gener-
 275 ated by GPT-4, into our training pipeline.

3.3. Architecture

281 **Model Design.** As illustrated in Fig. 4, our model consists
 282 of two main components: a multi-modal large language
 283

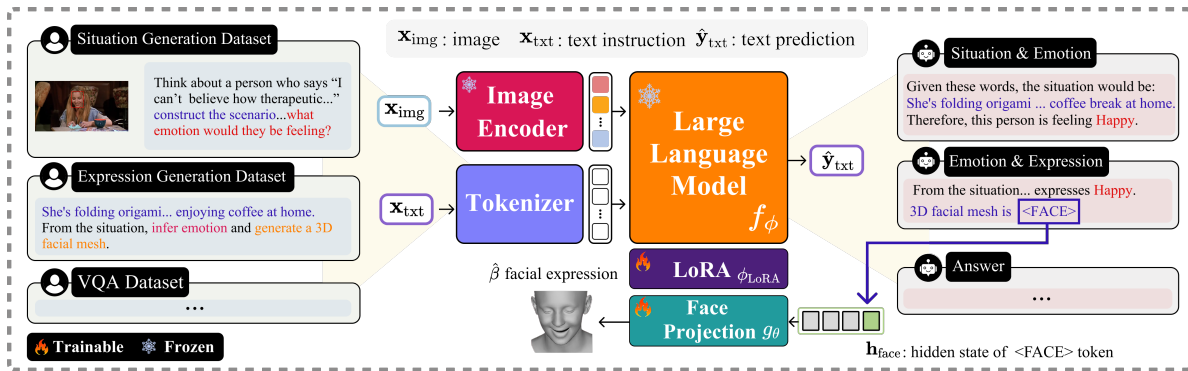


Figure 4. **Pipeline of ContextFace.** The instruction datasets described in Fig. 3 are used as input. Images and text are processed by the image encoder and tokenizer respectively, and are fed into the LLM to produce text output. For the expression generation task, the hidden state of the <FACE> token 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.

284 model denoted as f_ϕ and a face projection model g_θ . Text
 285 instructions \mathbf{x}_{txt} serve as inputs to the multi-modal LLM f_ϕ ,
 286 resulting in text outputs $\hat{\mathbf{y}}_{\text{txt}}$ given by

$$287 \quad \hat{\mathbf{y}}_{\text{txt}} = f_\phi(\mathbf{x}_{\text{txt}}). \quad (1)$$

288 We employed LoRA fine-tuning [14] with trainable param-
 289 eters ϕ_{LoRA} to efficiently adapt the model to our specific task
 290 while preserving the generation capabilities of pre-trained
 291 vision-language models. If the input prompt requests the
 292 generation of 3D facial mesh, the multi-modal LLM f_ϕ is
 293 trained to include <FACE> token in its text outputs $\hat{\mathbf{y}}_{\text{txt}}$.
 294 The hidden state of the multi-modal LLM f_ϕ corresponding
 295 to <FACE> token, denoted by \mathbf{h}_{face} , is used to generate the
 296 3D facial mesh. \mathbf{h}_{face} contains prompt-specific contextual
 297 information required to generate a facial expression. \mathbf{h}_{face}
 298 is subsequently passed through the face projection module
 299 g_θ to obtain the expression coefficient $\hat{\beta}$ given by

$$300 \quad \hat{\beta} = g_\theta(\mathbf{h}_{\text{face}}). \quad (2)$$

301 We use a two-layer multi-layer perceptron (MLP) for g_θ .
 302 **Loss Functions.** The loss function is designed to optimize
 303 both text output $\hat{\mathbf{y}}_{\text{txt}}$ and facial expression output $\hat{\beta}$. Firstly,
 304 the loss \mathcal{L}_{txt} for the text generation is the cross entropy loss
 305 between the predicted response $\hat{\mathbf{y}}_{\text{txt}}$ and the ground truth
 306 \mathbf{y}_{txt} :

$$307 \quad \mathcal{L}_{\text{txt}} = \text{CE}(\mathbf{y}_{\text{txt}}, \hat{\mathbf{y}}_{\text{txt}}). \quad (3)$$

308 Secondly, for the facial expression loss denoted by $\mathcal{L}_{\text{face}}$, we
 309 adopted the L2 loss as the objective function for our facial
 310 expression component. The loss compares L2 distance be-
 311 tween the predicted $\hat{\beta}$ and ground truth β coefficients given
 312 by

$$313 \quad \mathcal{L}_{\text{face}} = \|\beta - \hat{\beta}\|_2^2. \quad (4)$$

314 where $\|\cdot\|_2$ denotes the L2-norm of a vector. The overall
 315 loss \mathcal{L} consists of text generation loss \mathcal{L}_{txt} and facial expres-

sion $\mathcal{L}_{\text{face}}$, which is given by

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

where λ is a weight that balances the losses.

Training Process. ContextFace is jointly trained on three
 instruction datasets: Situation Generation, Expression Gen-
 eration, 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> to-
 ken. Also, the face projection g_θ is a trainable module
 to learn the mapping from \mathbf{h}_{face} to facial expression param-
 eters $\hat{\beta}$. Due to the inherently smaller magnitude of the facial
 expression loss values, we applied a weighting factor λ 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
 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 intro-
 duced 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 [29] and Qwen2.5-VL-7B [42]. ContextFace
 and all the baseline models are evaluated in a zero-shot
 manner on SFEW-C. In addition, we provide an evalua-
 tion on emotion datasets derived from MER2023 [4, 22]
 datasets: due to space constraints, readers are referred to
 Sec. 1 of Supplementary Materials.

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 [27]	96.50	<u>93.96</u>	<u>90.45</u>	77.78	73.33	78.05	71.43	81.42	84.98
LLaVA-NEXT-13B [29]	93.62	93.51	88.31	<u>91.57</u>	69.90	85.71	86.27	<u>86.88</u>	<u>89.73</u>
Qwen2.5-VL-7B [42]	<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	95.89	92.77	94.94	85.45	85.71	<u>84.31</u>	90.13	90.68
<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 [27]	<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 [29]	90.23	90.51	77.55	91.57	78.50	<u>83.72</u>	71.32	84.01	85.56
Qwen2.5-VL-7B [42]	90.65	97.96	91.72	<u>92.99</u>	80.70	76.36	<u>90.32</u>	<u>90.20</u>	<u>91.33</u>
ContextFace (ours)	96.60	92.21	<u>84.93</u>	94.48	91.74	95.65	92.78	93.21	94.89

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.

348 **Network Architecture.** We utilized LLaVA-v1.5-13B [27]
 349 as the backbone multimodal large language model for our
 350 implementation. We used the CLIP vision encoder [36]
 351 (ViT-L/14 with 336px resolution) of LLaVA-v1.5-13B for
 352 visual processing. The face projection module is a two-
 353 layer MLP with input, hidden layer and output sizes given
 354 by 5120, 5120, and 100 respectively.

355 4.2. Quantitative Results

356 **Contextual Emotion Recognition.** Tab. 1 shows the per-
 357 formance of contextual emotion recognition. Specifically,
 358 the task is to classify images accompanied by one of two
 359 context types, situation or quote, into seven categories of
 360 emotions. ContextFace achieves superior performance in
 361 the majority of emotion categories for both situation and
 362 quote contexts, with the highest UAR (Unweighted Average
 363 Recall) and WAR (Weighted Average Recall) scores in
 364 each contextual setting (90.13%, 90.68% with situation
 365 context and 93.21%, 94.89% with quote context). Notably,
 366 ContextFace excels at detecting emotions like Disgust and
 367 Fear which typical models for emotion recognition struggle
 368 to identify. This superior performance suggests that Con-
 369 textFace has developed a sophisticated understanding of the
 370 complex relationship between context and emotion through
 371 the proposed framework of visual instruction tuning.

372 **Facial Expression Prediction.** We assess the performance
 373 of predicting facial expression coefficient β for given situ-
 374 ation or context. The expression-context pairs in SFEW-C
 375 are used for evaluation. We measure the error in the pre-
 376 diction in terms of the distance between 100-dimensional
 377 coefficient vectors of facial expressions in the L2 and FD
 378 (Fréchet Distance) metrics. L2 measures the error in the
 379 L2 distance between the coefficient vectors. FD measures

Baseline	L2 ↓	FD ↓
<i>situation</i>		
Random	0.34 ± 0.0102	1.98 ± 0.2263
Mean	0.14 ± 0.0047	14.38 ± 0.4688
Ours	0.07 ± 0.0055	0.81 ± 0.0711
<i>quote</i>		
Random	0.31 ± 0.0093	2.42 ± 0.1982
Mean	0.12 ± 0.0042	11.76 ± 0.4178
Ours	0.1 ± 0.0048	1.95 ± 0.1305

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.

the statistical distance between the distributions of coefficient vectors. Since there exist no other methods for contextual expression generation, we perform a self-evaluation adopting the multiple axes-based evaluation framework discussed in previous studies [32, 33] 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

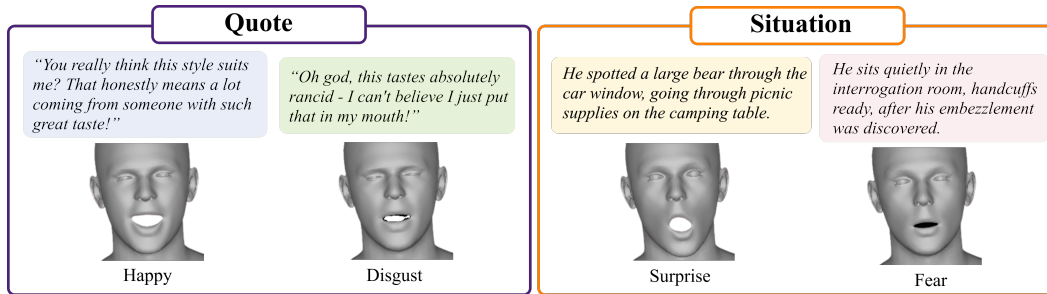


Figure 5. **Contextual Alignment of Emotional Expressions.** Generation results of facial expressions in response to quotes and situations.

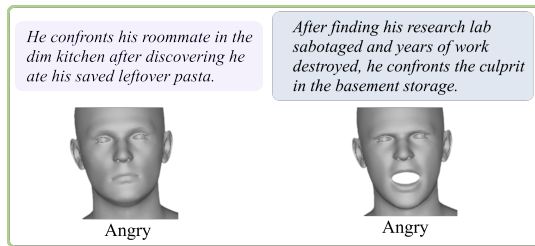


Figure 6. **Contextual Variation within Emotions.** Contextual variations of angry facial expressions in different situations.

Situation Generation	Expression Generation	UAR	WAR
×	×	81.42	84.98
✓	×	79.34	82.48
×	✓	88.33	89.40
✓	✓	93.21	94.89

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

393 superior FD scores compared to the Mean baseline, preserv-
 394 ing the distinctive elements of specific facial expressions to
 395 effectively capture fine facial feature details.

396 4.3. Qualitative Results

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

401 **Intra-Emotion Expressional Diversity.** In Fig. 6, we ob-
 402 serve that ContextFace is capable of generating different
 403 expressions within the same emotion depending on the sit-
 404 uational context. Facial expressions involving anger were
 405 generated for different situations. The left image displays a
 406 milder expression in response to a minor annoyance, while
 407 the right image shows a significantly more intense expres-
 408 sion of rage in response to an infuriating situation.

409 **Speculative Expression Generation through Emotional**
 410 **Reasoning.** We show that our model can generate logically
 411 valid situations when given a quote for ambiguous facial
 412 expressions. Additionally, based on these generated situa-
 413 tions, the model can create appropriate facial expressions.
 414 In Fig. 7, we observe that our model correctly infers the
 415 emotions of the speaker, even with the same image and dif-
 416 ferent quotes. Moreover, the model can infer the emotional
 417 state or facial expression of not only the speaker, but also
 418 the other person. As shown in the example, while the quote
 419 comes from a woman, the model can *speculatively* generate
 420 the facial expression of her boyfriend through an interactive

prompting. This suggests that our model has emotional rea- 421
 soning capabilities to understand interpersonal relationships 422
 and infer corresponding emotions or facial expressions. 423

Facial expression generation for occluded faces. As 424
 shown in Fig. 8, our model can infer the situation based 425
 on contexts (image and quote), and can generate plausible 426
 facial expressions, even when the face is intentionally and 427
 completely occluded. Our model successfully infers situa- 428
 tions and generates appropriate expressions, demonstrat- 429
 ing it has comprehensively learned the relationships among 430
 context, emotions, and facial expressions. 431

432 4.4. Ablation study

Impact of Dataset Combinations. As shown in Fig. 3, 433
 we used three datasets to jointly train ContextFace, i.e., two 434
 context-augmented datasets: Situation Generation Dataset, 435
 Expression Generation Dataset and one VQA Dataset. We 436
 measured the model’s ability to infer emotions with or with- 437
 out context-augmented datasets where VQA dataset is used 438
 by default. As shown in Tab. 3, the model performed the 439
 best when all three datasets are used. The model without 440
 Situation Generation dataset performs marginally below the 441
 baseline LLaVA-1.5, indicating that balanced training with 442
 generalizable datasets rather than task-specific ones is es- 443
 sential for maintaining emotion inference capabilities. 444

Impact of Model Scale. The LLaVA-1.5-13B model 445
 demonstrated higher performance in both UAR and WAR 446
 compared to LLaVA-1.5-7B, as shown in Tab. 4. This sug- 447
 gests that as model size increases, emotional reasoning and 448
 expression generation capabilities improve. However, the 449

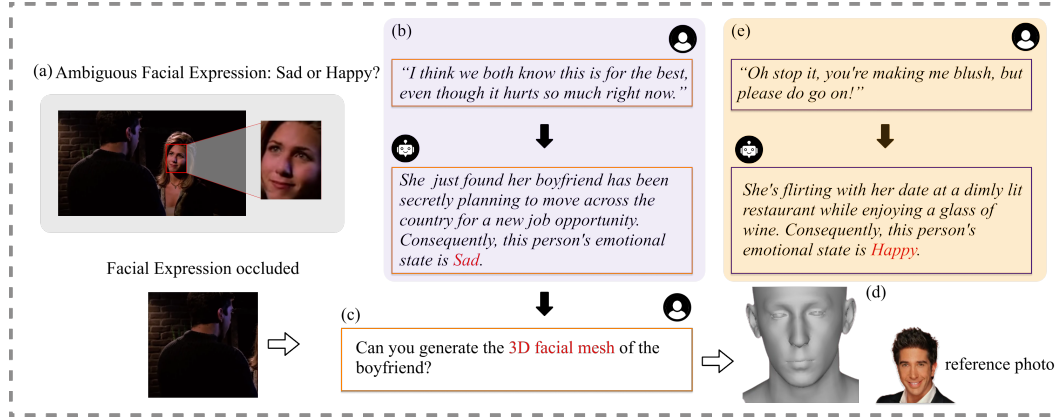


Figure 7. **Speculative Expression Generation through Emotional Reasoning.** (a) The expression in the input image is ambiguous. (b) Given the woman’s quote, our model can generate a plausible situation from the image along with the proper emotion. (c) The model is prompted to **speculate** the expression 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.

Method	UAR (%)	WAR
LLaVA-1.5-7B	92.00	91.70
LLaVA-1.5-13B	93.21	94.89

Table 4. **Ablation study.** Comparing emotion inference performance between LLaVA-1.5-13B and LLaVA-1.5-7B.

450 7B model maintains performance above 90% of the 13B
451 model, making it a computationally efficient alternative.

452 **5. Conclusion**

453 In this paper, we proposed ContextFace, a model that un-
454 derstands the deep association between contextual emotions
455 and facial expressions and generates context-aligned facial
456 expressions. Furthermore, we introduce new CAER-S-C
457 and SFEW-C datasets that enable the proposed contextual
458 learning process. To our knowledge, we are the first to
459 extend MLLMs beyond textual emotion analysis to direct

generation of facial expressions, bridging the gap between 460
language-driven emotion interpretation and visual synthesis. 461
The proposed integration will become increasingly 462
valuable as AI services evolve towards human-like inter- 463
actions through emotionally aware and visually expressive 464
communication. 465

466 **References**

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, An- 467
toine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur 468
Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, 469
Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, 470
Sina Samangooei, Marianne Monteiro, Jacob Menick, Se- 471
bastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sa- 472
hand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, 473
Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 474
Flamingo: a visual language model for few-shot learning, 475
2022. 2 476

[2] Anthropic. Claude-3.5-sonnet-20241022. [https://www. 477
anthropic.com](https://www.anthropic.com), 2024. Large Language Model API. 3 478

- 479 [3] Zehranaz Canfes, M Furkan Atasoy, Alara Dirik, and Pinar
480 Yanardag. Text and image guided 3d avatar generation and
481 manipulation. In *Proceedings of the IEEE/CVF Winter Con-*
482 *ference on Applications of Computer Vision*, pages 4421–
483 4431, 2023. 2
- 484 [4] Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Jingdong Sun,
485 Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and
486 Alexander Hauptmann. Emotion-llama: Multimodal emo-
487 tion recognition and reasoning with instruction tuning, 2024.
488 1, 2, 5
- 489 [5] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat
490 Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale
491 Fung, and Steven Hoi. Instructblip: Towards general-
492 purpose vision-language models with instruction tuning,
493 2023. 5, 6
- 494 [6] Radek Danecek, Michael J. Black, and Timo Bolkart.
495 EMOCA: Emotion driven monocular face capture and an-
496 imation. In *Conference on Computer Vision and Pattern*
497 *Recognition (CVPR)*, pages 20311–20322, 2022. 2
- 498 [7] Radek Daněček, Timo Bolkart, and Wojciech Zielonka. In-
499 ferno. <https://github.com/radekd91/inferno>,
500 2020. 1, 2, 4, 8
- 501 [8] Radek Daněček, Kiran Chhatre, Shashank Tripathi, Yan-
502 dong Wen, Michael Black, and Timo Bolkart. Emotional
503 speech-driven animation with content-emotion disentangle-
504 ment. In *SIGGRAPH Asia 2023 Conference Papers*, page
505 1–13. ACM, 2023. 2
- 506 [9] Abhinav Dhall, Roland Goecke, Simon Lucey, and Tom
507 Gedeon. Static facial expression analysis in tough condi-
508 tions: Data, evaluation protocol and benchmark. In *2011*
509 *IEEE International Conference on Computer Vision Work-*
510 *shops (ICCV Workshops)*, pages 2106–2112, 2011. 2, 3, 5
- 511 [10] Yasaman Etesam, Özge Nilay Yalçın, Chuxuan Zhang, and
512 Angelica Lim. Contextual emotion recognition using large
513 vision language models. In *2024 IEEE/RSJ International*
514 *Conference on Intelligent Robots and Systems (IROS)*, pages
515 4769–4776. IEEE, 2024. 2
- 516 [11] Yasaman Etesam, Özge Nilay Yalçın, Chuxuan Zhang, and
517 Angelica Lim. Contextual emotion recognition using large
518 vision language models, 2024. 1
- 519 [12] Yao Feng, Jing Lin, Sai Kumar Dwivedi, Yu Sun, Priyanka
520 Patel, and Michael J. Black. ChatPose: Chatting about 3d
521 human pose. In *CVPR*, 2024. 2, 3
- 522 [13] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar
523 Hilliges. Vid2avatar: 3d avatar reconstruction from videos
524 in the wild via self-supervised scene decomposition. In *Pro-*
525 *ceedings of the IEEE/CVF Conference on Computer Vision*
526 *and Pattern Recognition*, pages 12858–12868, 2023. 2
- 527 [14] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-
528 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
529 Lora: Low-rank adaptation of large language models, 2021.
530 5
- 531 [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui
532 Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation
533 via large language model. *arXiv preprint arXiv:2308.00692*,
534 2023. 2, 3
- [16] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and
Kwanghoonn Sohn. Context-aware emotion recognition net-
works. In *Proceedings of the IEEE/CVF international con-*
ference on computer vision, 2019. 2, 3, 5
- [17] Yuxuan Lei, Dingkang Yang, Zhaoyu Chen, Jiawei Chen,
Peng Zhai, and Lihua Zhang. Large vision-language mod-
els as emotion recognizers in context awareness, 2024. 1
- [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 ex-
pression 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 Con-*
ference on Learning Representations. 1
- [21] Yifan Li, Anh Dao, Wentao Bao, Zhen Tan, Tianlong Chen,
Huan Liu, and Yu Kong. Facial affective behavior analysis
with instruction tuning, 2024. 1
- [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. 5
- [23] Zheng Lian, Licai Sun, Mingyu Xu, Haiyang Sun, Ke Xu,
Zhuofan Wen, Shun Chen, Bin Liu, and Jianhua Tao. Ex-
plainable multimodal emotion reasoning. *CoRR*, 2023. 2
- [24] Zheng Lian, Haiyang Sun, Licai Sun, Jiangyan Yi, Bin Liu,
and Jianhua Tao. Affectgpt: Dataset and framework for ex-
plainable multimodal emotion recognition, 2024. 1
- [25] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Moham-
mad Shoeybi, and Song Han. Vila: On pre-training for vi-
sual language models. In *Proceedings of the IEEE/CVF con-*
ference on computer vision and pattern recognition, pages
26689–26699, 2024. 1
- [26] Yihong Lin, Liang Peng, Jianqiao Hu, Xiandong Li, Wenx-
iong 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, Qingyang Wu, and Yong Jae Lee.
Visual instruction tuning. In *NeurIPS*, 2023. 1, 2, 3, 4, 5, 6
- [28] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
Improved baselines with visual instruction tuning. In *Pro-*
ceedings of the IEEE/CVF Conference on Computer Vision
and Pattern Recognition, pages 26296–26306, 2024. 2
- [29] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan
Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Im-
proved reasoning, ocr, and world knowledge, 2024. 2, 5,
6
- [30] 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. 1

- 592 [31] Zhentao Liu, Min Wu, Weihua Cao, Luefeng Chen, Jianping
593 Xu, Ri Zhang, Mengtian Zhou, and Junwei Mao. A facial ex-
594 pression emotion recognition based human-robot interaction
595 system. *IEEE CAA J. Autom. Sinica*, 4(4):668–676, 2017. 1
- 596 [32] Evonne Ng, Hanbyul Joo, Liwen Hu, Hao Li, Trevor Darrell,
597 Angjoo Kanazawa, and Shiry Ginosar. Learning to listen:
598 Modeling non-deterministic dyadic facial motion, 2022. 6
- 599 [33] Evonne Ng, Sanjay Subramanian, Dan Klein, Angjoo
600 Kanazawa, Trevor Darrell, and Shiry Ginosar. Can language
601 models learn to listen?, 2023. 6
- 602 [34] Mang Ning, Albert Ali Salah, and Itir Onal Ertugrul. Repre-
603 sentation learning and identity adversarial training for facial
604 behavior understanding, 2024. 1
- 605 [35] Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu
606 Zhu, Hongyan Liu, Jun He, and Zhaoxin Fan. Emotalk:
607 Speech-driven emotional disentanglement for 3d face anima-
608 tion. In *Proceedings of the IEEE/CVF international confer-
609 ence on computer vision*, 2023. 2
- 610 [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
611 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
612 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen
613 Krueger, and Ilya Sutskever. Learning transferable visual
614 models from natural language supervision, 2021. 6
- 615 [37] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdel-
616 rahman Shaker, Salman Khan, Hisham Cholakkal, Rao M.
617 Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S. Khan.
618 Glamm: Pixel grounding large multimodal model. *The
619 IEEE/CVF Conference on Computer Vision and Pattern
620 Recognition*, 2024. 2, 3
- 621 [38] Niyati Rawal and Ruth Maria Stock-Homburg. Facial emo-
622 tion expressions in human–robot interaction: A survey. *In-
623 ternational Journal of Social Robotics*, 14(7):1583–1604,
624 2022. 1
- 625 [39] Arnab Kumar Roy, Hemant Kumar Kathania, Adhitiya
626 Sharma, Abhishek Dey, and Md. Sarfaraj Alam Ansari. Re-
627 semotenet: Bridging accuracy and loss reduction in facial
628 emotion recognition. *IEEE Signal Processing Letters*, pages
629 1–5, 2024. 1
- 630 [40] Matteo Spezialetti, Giuseppe Placidi, and Silvia Rossi. Emo-
631 tion recognition for human-robot interaction: Recent ad-
632 vances and future perspectives. *Frontiers in Robotics and
633 AI*, 7:532279, 2020. 1
- 634 [41] Andreas Steiner, André Susano Pinto, Michael Tschannen,
635 Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Grit-
636 senko, Matthias Minderer, Anthony Sherbondy, Shangbang
637 Long, et al. Paligemma 2: A family of versatile vlms for
638 transfer. *arXiv preprint arXiv:2412.03555*, 2024. 2
- 639 [42] Qwen Team. Qwen2.5-vl, 2025. 2, 5, 6
- 640 [43] Chengpeng Wang, Li Chen, Lili Wang, Zhaofan Li, and Xue-
641 bin Lv. Qcs: Feature refining from quadruplet cross similar-
642 ity for facial expression recognition, 2025. 1
- 643 [44] Alexandros Xenos, Niki Maria Foteinopoulou, Ioanna Nti-
644 nou, Ioannis Patras, and Georgios Tzimiropoulos. Vllms pro-
645 vide better context for emotion understanding through com-
646 mon sense reasoning, 2024. 1, 2
- 647 [45] Bohao Xing, Zitong Yu, Xin Liu, Kaishen Yuan, Qilang
648 Ye, Weicheng Xie, Huanjing Yue, Jingyu Yang, and Heikki
Kälviäinen. Emo-llama: Enhancing facial emotion under-
standing with instruction tuning, 2024. 1
- [46] Yazhou Zhang, Mengyao Wang, Youxi Wu, Prayag Tiwari,
Qiuchi Li, Benyou Wang, and Jing Qin. DialogueUlm: Con-
text and emotion knowledge-tuned large language models for
emotion recognition in conversations, 2024. 1, 2