

Evaluating Web-trained Facial Expression Recognition in Collaborative Problem-Solving

Anonymous CVPR submission

Paper ID 21

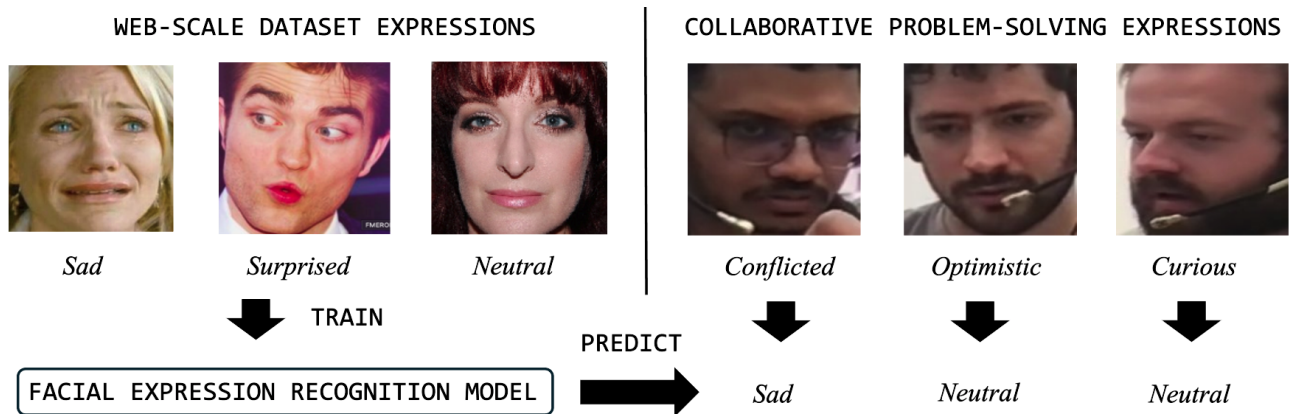


Figure 1. Popular facial expression recognition (FER) models are trained to predict canonical basic emotions from prototypical facial expression images, often obtained from web queries and annotated with external observation or query-matching. In collaborative problem-solving (CPS), epistemic emotions such as *curiosity*, *optimism*, and *conflict* are often more relevant but manifest more subtly in facial behavior and are subjective in their nature. We examine the alignment between these model outputs and epistemic affect observed in collaborative problem-solving settings. Samples from the left group are from AffectNet [1] and from the right are from our CPS dataset.

Abstract

001 Facial expression recognition (FER) have been adopted ex-
 002 tensively in educational research. These models are usu-
 003 ally trained on web-scale datasets optimized for cano-
 004 nical basic emotions. However, in collaborative learn-
 005 ing, the affective states of interest are often epistemic
 006 (learning-relevant) in nature—such as confusion, curios-
 007 ity, and frustration—rather than prototypical basic emo-
 008 tions like disgust or anger. Here, we evaluate popular
 009 FER systems trained using web-scale datasets on small-
 010 group collaborative problem-solving sessions. We anno-
 011 tate epistemic emotions using a retrospective cued-recall
 012 approach, where participants watch a video of a collabo-
 013 rative session immediately after the group task and individ-
 014 ually self-report at different segments. We analyze cross-
 015 taxonomy alignment between epistemic emotions and pre-
 016 dicted basic emotions, cross-model agreement in collab-
 017 orative learning, and dimensional valence–arousal struc-

ture with respect to basic and epistemic emotions. Across 018
 models, categorical predictions overwhelmingly collapse 019
 to Neutral for all epistemic labels, and valence–arousal 020
 outputs fail to distinguish the different epistemic states. 021
 Furthermore, low agreement between models indicates 022
 some degree of instability in both web-scale and natu- 023
 ralist collaborative contexts. Our findings suggest that 024
 FER systems trained for canonical emotion recognition on 025
 web-scale datasets do not directly generalize to the more 026
 subtle epistemic emotions experienced during collabora- 027
 tion, which has implications for education researchers de- 028
 ploying these tools in classrooms. Our code is avail- 029
 able at <https://anonymous.4open.science/r/cvpr-edu-affect-2026-6EF5>. 030
 031

1. Introduction

Recent advances in computer vision have increased the 033
 adoption of facial expression recognition (FER) tools in 034

035 education. Open-source frameworks [2, 3], large-scale
036 datasets [1, 4], and modern architectures [5] have lowered
037 the barrier to deploying FER systems in applications such
038 as classroom sensing and intelligent tutoring systems [6–9].
039 Despite this rapid adoption, relatively little empirical work
040 has examined how web-trained FER systems perform in
041 naturalistic multi-party learning environments such as col-
042 laborative problem-solving (CPS). Most FER models are
043 trained to predict canonical basic emotions—such as hap-
044 piness, sadness, anger, fear, disgust, and surprise—using
045 datasets curated primarily from scraping or querying the
046 web [1, 10] and annotated by external coders with subsets
047 including one annotator per image [11] or computer vision
048 algorithms that are trained on similar manual annotations
049 [4]. Facial expressions from these datasets tend to be more
050 prototypical and exaggerated compared to naturalistic col-
051 laborative data since many are stock images of posed ex-
052 pressions or celebrities at televised events, and have been
053 criticized for being inconsistent in their labels due to their
054 unreliable annotation process [11]. Recent affective com-
055 puting research has also expanded to dimensional repre-
056 sentations, most notably the valence–arousal (V–A) cir-
057 cumplex model [12]. Many modern FER datasets provide
058 both categorical and continuous annotations [1, 11], en-
059 abling models to predict valence and arousal alongside dis-
060 crete emotions. Such datasets typically also include soft-
061 labels (the probability vectors for discrete labels) with each
062 instance, and offers a potentially richer representational
063 space [13].

064 In contrast, the affective states more relevant to educa-
065 tional contexts such as CPS are epistemic in nature, includ-
066 ing confusion, curiosity, frustration, and engagement, and
067 they manifest in facial expression much more subtly. These
068 emotions are closely tied to knowledge construction, uncer-
069 tainty management, and collaborative coordination, and
070 have been shown to influence learning outcomes in various
071 educational contexts [14–16]. Consequentially, they are
072 of particular interest for instructors and learning technol-
073 ogies that aim to monitor and support learning processes in
074 real time [14]. To that end, researchers are using machine
075 learning with behavioral data that can unobtrusively cap-
076 ture these states during authentic learning activities [17].
077 FER systems are becoming promising candidates for this
078 purpose due to their increasing efficiency and ease of adop-
079 tion [6].

080 These trends bring some important questions to atten-
081 tion: when web-trained FER models are used in collab-
082 orative learning settings, do they meaningfully align with
083 epistemic affective states or collapse toward dominant ba-
084 sic emotion categories? Moreover, do the continuous va-
085 lence–arousal outputs of VA models capture finer affective
086 distinctions in this context, or do they largely reflect the
087 same categorical structure embedded in their training data?

088 Finally, do different FER systems/architectures even pro-
089 duce consistent signals under such a domain shift?

090 To address these questions, we evaluate multiple widely
091 used pretrained FER systems on a dataset of CPS ses-
092 sions annotated with self-reported epistemic states collected
093 via retrospective cued recall [18, 19]. We analyze cross-
094 taxonomy alignment, prediction concentration and inter-
095 model agreement. Our goal is to examine how contempo-
096 rary FER systems behave out-of-the-box when deployed in
097 CPS settings.

098 2. Related Work

099 The last decade has seen significant interest in FER, primar-
100 ily driven by web-scale datasets and deep learning. Early
101 large-scale annotated datasets such as FER2013 [10] en-
102 abled supervised training of CNNs for basic emotion recog-
103 nition. AffectNet [1] and RAF-DB [20] further expanded
104 the scale of training data and introduced both categori-
105 cal and valence–arousal annotations, with AffectNet+ [11]
106 also introducing soft labels for facial expressions. Re-
107 cent architectures have even leveraged attention mecha-
108 nisms and transformer-style backbones to improve perfor-
109 mance on these benchmarks. For example, dual-branch at-
110 tention networks [21] and ViT based architectures such as
111 POSTER++ [5] reported strong results on multiple in-the-
112 wild emotion datasets. Concurrently, toolkits such as Open-
113 Face [2, 3] and LibreFace [22] have been developed by re-
114 searchers to provide user-friendly and efficient interfaces
115 for end-to-end extraction of facial landmarks, action units,
116 gaze, head pose estimation and facial expression prediction,
117 facilitating adoption in domains such as education.

118 Most training and evaluation of FER models remain
119 within the distribution of the web-scale benchmark datasets.
120 Such datasets are often constructed by querying the web
121 with emotion-related keywords [1, 4, 10] using tools such
122 as WordNet [23]. This process have yielded posed or ex-
123 aggerated expressions that approximate canonical emotion
124 prototypes [13]. Such data differ substantially from sponta-
125 neous affect observed in real-world interaction in terms of
126 intensity, dynamics, and contextual embedding. In collabo-
127 rative learning settings, learners often regulate or mask vis-
128 ible affect during group interaction [24]. Negative epistemic
129 states such as confusion or frustration may manifest subtly
130 or co-occur with task-focused expressions. Group dynamics
131 introduce additional layers of co-regulation and social mon-
132 itoring [25]. As a result, the mapping between visible facial
133 configurations and epistemic emotions may differ from that
134 observed in web-scale datasets. In the learning sciences, af-
135 fect is often studied not as isolated emotional expressions
136 but as part of a broader cognitive–affective process that un-
137 folds during knowledge construction. Epistemic emotions
138 such as confusion, curiosity, and frustration are closely tied
139 to processes such as impasse resolution, hypothesis gen-

140 eration, and collaborative sense-making [14, 15]. These
141 states are frequently transient, low-intensity, and embedded
142 within ongoing task activity, making them difficult to infer
143 from facial cues alone.

144 Many web-scaled datasets also use the basic emotion
145 taxonomies and/or the dimensional affect models to repre-
146 sent facial affect. Basic emotion theory characterizes af-
147 fect using a small set of discrete categories—such as hap-
148 piness, sadness, anger, fear, disgust, and surprise—that are
149 assumed to correspond to recognizable facial configurations
150 and are widely used in FER datasets and benchmarks [10,
151 26]. Datasets using this taxonomy often use a “Neutral”
152 label if the annotator believes none of the basic emotions
153 are present. Dimensional approaches, most notably the va-
154 lence–arousal circumplex model, instead represent affect
155 along continuous axes capturing positivity–negativity (va-
156 lence) and activation level (arousal) [12]. These represen-
157 tations have enabled FER systems to predict both categori-
158 cal and continuous affective states from facial images [1].
159 In educational contexts, it remains unclear whether the
160 epistemic states of interest correspond directly to canon-
161 ical facial prototypes associated with basic emotions, or
162 map cleanly onto simple regions of valence–arousal space.
163 Understanding how existing FER representations relate to
164 these learning-relevant affective states is therefore an im-
165 portant step to consider before deploying such models in
166 collaborative learning settings.

167 Another concerning factor is the annotation of affect in
168 educational contexts. Many FER benchmarks rely on an-
169 notations from external observers who infer emotional states
170 from static images or short clips [1, 20]. Although inter-
171 rater agreement is often reported, observer-based labeling
172 may not fully capture internal experiences, particularly in
173 complex social contexts. Some subsets of popular datasets
174 also only have one annotator assigned per image due to the
175 sheer volume of images. Benitez-Quiroz et al. [4] have
176 tried to improve this paradigm by generating annotations
177 for millions of facial images using automated facial be-
178 havior analysis systems trained on manually coded Facial
179 Action Units. While such approaches enable dataset scal-
180 ing, subsequent work has raised concerns about the reliabil-
181 ity of labels generated by automated FER systems, noting
182 that algorithmic annotations may propagate biases and sys-
183 tematic errors present in the underlying models [27]. Ed-
184 ucational research has therefore explored alternative an-
185 notation strategies, including self-reports, behavioral coding,
186 and multimodal sensing approaches [7, 14]. Physiological
187 signals such as electrodermal activity [28] and heart rate
188 variability [29] have been used to study learning-related af-
189 fect, but these methods are often intrusive and impractical
190 in authentic learning environments. A recently adopted ap-
191 proach is retrospective cued-recall (RCR), in which partic-
192 ipants review recordings of their own activity and report

Table 1. Demographic information of participants in the EECPS-WT dataset [19]

Gender	Male (18)	Female (7)	Non-binary (2)
Age	18–24 (14)	25–34 (11)	35+ (2)
Ethnicity	Asian (13)	White (12)	Hispanic (2)



Figure 2. Three team members collaborating during an experimen-
tial task from the EECPS-WT dataset

193 their internal states at specific moments in time [18]. This
194 method preserves contextual grounding while allowing partic-
195 ipants to reflect on their affective experiences during col-
196 laboration. Such annotations provide a complementary per-
197 spective to observer-based labels and offer a practical mech-
198 anism for studying epistemic emotions in naturalistic set-
199 tings.

3. Methods 200

3.1. Dataset 201

202 We use a collaborative learning dataset, Epistemic Emo-
203 tions in Collaborative Problem Solving: Weights Task
204 (EECPS-WT). This dataset is a contribution from previ-
205 ous work [redacted for review] which we build upon in
206 this paper. The data is collected from 27 college students
207 who were over 18 years old and fluent in English. Table
208 1 shows their demographic information. Participants
209 were organized into nine groups of three and were video
210 recorded while completing a collaborative problem-solving
211 (CPS) activity, adapted from Khebour et al [30]. In this
212 activity, groups used a balance scale to infer the relative
213 weights of five blocks, one of which has a known reference
214 weight. Each group was required to submit a single shared
215 answer via survey upon completion of the phase. Interac-
216 tions were recorded using consumer-grade laptop webcams
217 positioned in front of participants, resulting in naturalistic,
218 low-resource video data similar to what may be encountered
219 in classroom or remote learning settings. Figure 2 shows an
220 example frame from the dataset, where three team members
221 are collaborating on the task.

222 Following task completion, participants individually
 223 viewed a recording of their collaborative session using a ret-
 224 rospective cued recall (RCR) procedure. During playback,
 225 participants reported their own affective states using an in-
 226 teractive survey interface. Reports could be self-initiated at
 227 any time. Additionally, if no report was submitted for 60
 228 seconds, the system automatically issued a probe request-
 229 ing a report. This design captures both spontaneous and
 230 probe-caught affective states. Participants selected one or
 231 more affective states from a predefined set of seven epis-
 232 temic labels: Confused, Curious, Frustrated, Disengaged,
 233 Optimistic, Surprised, and Conflicted, adopted from prior
 234 research [31]. Each selected label was timestamped accord-
 235 ing to video playback time. All reports were standardized
 236 and consolidated into a single dataset, with each affective
 237 label represented as a separate instance. Reports were ad-
 238 ditionally tagged as self-caught or probe-caught based on
 239 their trigger mechanism. Statistical information such as fre-
 240 quency distribution and temporal distribution of labels is re-
 241 ported in Sec 1 of the supplementary material.

242 3.2. Data Preprocessing

243 To evaluate facial expression recognition (FER) models at
 244 the level of individual reports, we construct face-based im-
 245 age instances aligned with self-reported timestamps. All
 246 videos were processed frame-by-frame using RetinaFace
 247 for face detection [32]. For each detected face, a bound-
 248 ing box was extracted and cropped from the original frame.
 249 Since participants were seated in fixed positions during
 250 the collaborative task, with participants 1 through 3 seated
 251 from left to right from the camera’s point of view, de-
 252 tectations were organized per participant by spatial consis-
 253 tency across frames using the x-coordinate of detections
 254 to map them to participants. Cropped faces were stored
 255 at one-second resolution, with each second containing ap-
 256 proximately 30 frames corresponding to the recorded frame
 257 rate. All subsequent model inference operates exclusively
 258 on these cropped face images.

259 Each report is associated with a specific playback times-
 260 tamp. To account for potential temporal misalignment be-
 261 tween when an affective state was experienced and when
 262 it was reported, we adopt a temporal neighborhood sam-
 263 pling strategy. For each affect report at time t , we collect all
 264 cropped face images within a ± 5 second window, i.e., the
 265 interval $[t-5, t+5]$. From this pool of candidate frames, we
 266 randomly sample 10 face images to represent the instance.
 267 This strategy reduces sensitivity to single-frame noise and
 268 accommodates slight reporting delays inherent in retrospec-
 269 tive annotation. Each affect report is therefore represented
 270 as a set of 10 independently processed face images. Models
 271 are applied independently to each of the 10 frames, and pre-
 272 dictions are aggregated to produce an instance-level output
 273 (described in Section 3.3).

274 3.3. Facial Expression Recognition Models

275 We benchmark FER models that output (i) categorical ba-
 276 sic emotions or (ii) continuous valence–arousal representa-
 277 tions. We perform inference independently on each sampled
 278 face crop, and predictions are aggregated at the instance
 279 level. For a given affect report at time t , let $\{x_i\}_{i=1}^K$ de-
 280 note the set of $K = 10$ sampled face images within the
 281 temporal window described in Section 3.2. Each image is
 282 resized to 224×224 pixels and normalized according to the
 283 preprocessing requirements of the corresponding model.

284 3.3.1. Categorical Basic Emotion Models

285 We evaluate multiple pretrained FER models that predict
 286 discrete basic emotion categories. These include both
 287 multi-task pipelines OpenFace 3.0 [3] and LibreFace [22]
 288 and off-the-shelf deep learning FER models POSTER++
 289 [5], EmotiEffLib [33], and DDAMFN [21]. In addition, we
 290 include Qwen2.5-VL-7B [34], a vision language model, as
 291 a general purpose comparison point. OpenFace-3.0 outputs
 292 class probabilities (also known as soft labels [11]) and the
 293 other five outputs hard labels. All six models are used with
 294 publicly released pretrained weights.

295 For models that output soft labels, each image x_i pro-
 296 duces a probability vector:

$$297 \mathbf{p}_i = f(x_i) \in \mathbb{R}^C, \quad \text{with} \quad \sum_{c=1}^C p_{i,c} = 1.$$

298 To obtain an instance-level prediction, we compute the
 299 mean probability vector across sampled frames:

$$300 \bar{\mathbf{p}} = \frac{1}{K} \sum_{i=1}^K \mathbf{p}_i,$$

301 and assign the predicted label as

$$302 \hat{y} = \arg \max_c \bar{p}_c.$$

303 Averaging probabilities preserves confidence informa-
 304 tion across frames and reduces sensitivity to single-frame
 305 noise.

306 For models that return only a discrete label per image,
 307 we aggregate predictions via majority vote:

$$308 \hat{y} = \text{mode}\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_K\}.$$

309 In the case of ties, a deterministic class ordering is used.
 310 Majority voting is appropriate when per-frame confidence
 311 scores are not available or not directly comparable across
 312 models.

313 For Qwen 2.5-VL, we provide the 10 frames along with
 314 this text prompt: *These images are frames from a single*
 315 *video. Classify the emotion being displayed. Choose ex-*
 316 *actly one from the universal emotions. Output ONLY the*
 317 *emotion name, and nothing else.*

Because emotion labels vary slightly across different model outputs (e.g., “Happiness” vs. “Happy”), predicted categories are mapped to a standardized set: $\{Neutral, Happy, Sad, Angry, Surprise, Fear, Disgust, Contempt\}$. *Contempt* was removed during cross-model analysis since some categorical models does not use it as one of their outputs [35]. This normalization enables cross-model comparison under a unified taxonomy.

3.3.2. Dimensional Valence–Arousal Model

To examine whether continuous affect representations better align with collaborative learning contexts, we evaluate a pretrained system that jointly predicts discrete emotions and continuous valence–arousal values: EmotiEffLib [33]. For each image x_i , the model outputs:

$$(v_i, a_i) = f_{VA}(x_i),$$

where $v_i \in [-1, 1]$ denotes valence and $a_i \in [-1, 1]$ denotes arousal.

Instance-level dimensional predictions are computed by averaging across sampled frames:

$$(\bar{v}, \bar{a}) = \frac{1}{K} \sum_{i=1}^K (v_i, a_i).$$

This produces a single valence–arousal coordinate for each affect report instance.

We evaluate a single representative dimensional model rather than benchmarking multiple architectures. Because valence–arousal models are typically trained on the same web-scale datasets [1, 36] and share similar objective functions, we picked one model as our representative to examine whether continuous affect representations meaningfully differentiate epistemic states in collaborative learning.

3.4. Analysis

Our experiments focus on understanding how FER models perform on collaborative learning data. Because epistemic emotions and basic emotions/valence-arousal represent different conceptual spaces, we do not treat this as a conventional supervised classification problem. Instead, we analyze model behavior through distributional alignment, prediction concentration, and cross-model agreement.

3.4.1. Categorical Affect Alignment

We first examine how predicted basic emotion categories distribute across epistemic labels. For each categorical model, we compute a cross-taxonomy confusion matrix whose rows correspond to epistemic emotions (e.g., Curious, Confused, Frustrated) and columns correspond to predicted basic emotions. To account for the different distributions of labels in the dataset, we report the row-normalized version of this matrix, representing the conditional distribution of predicted basic emotions given each epistemic state.

This analysis allows us to assess whether specific epistemic states map consistently onto particular basic emotions, or whether predictions collapse toward dominant categories. To maintain clarity and conciseness, we only report the row-normalized confusion matrix for OpenFace 3.0 as the representative categorical model and summarize results across all other evaluated models in Sec 2.1 of the supplementary material. We also compute the proportion of instances assigned to each basic emotion category and report the most frequent predicted basic emotion for each epistemic label across all categorical models.

3.4.2. Dimensional Affect Structure

We next analyze continuous valence–arousal predictions produced by the dimensional model. For each epistemic label, we compute the mean and standard deviation of valence and arousal values and visualize their distributions using scatter plots in the two-dimensional affect space. This analysis addresses whether dimensional representations provide separable structure for epistemic states that categorical models fail to capture. We examine if epistemic states form clusters in valence–arousal space, the degree of overlap between states, and the stability of dimensional predictions across instances. In addition, we analyze valence–arousal distributions grouped by predicted basic emotion categories to verify internal coherence of the dimensional outputs and how the basic emotions complement the dimensional predictions. Here, we report scatter plots in the valence-arousal space to represent the distributions from both epistemic labels and basic emotion labels. Detailed statistical reports of these distributions and box plots are in Sec 2.2 of the supplementary material.

3.4.3. Cross-Model Agreement

To assess robustness across architectures under domain shift, we computed pairwise confusion matrices between categorical predictions generated by two FER models on the same instances. This analysis was conducted across random samples from EECPS-WT and AffectNet [1], a large-scale, in-the-wild facial expression benchmark. Pairwise confusion matrices were then constructed by treating the Model 1 predictions as the reference axis and the Model 2 predictions as the comparison axis. We normalize the counts of predictions from Model 1 for each label to counter the influence of label imbalance in evaluating cross-model agreement. This allows us to directly compare inter-model agreement between these FER models in two different datasets, one that is web-scaled and the other educational. Here, we report confusion matrices representing agreement between OpenFace 3.0 and LibreFace as the representative models. Confusion matrices between all other combinations of models is reported in Sec 2.3 of the supplementary material.

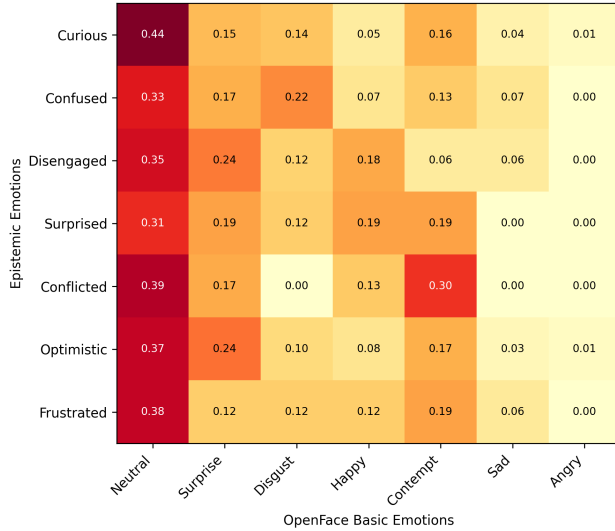


Figure 3. Row-normalized confusion matrix between epistemic states and OpenFace 3.0 predictions.

4. Results

We report results for our experiments in: (1) cross-taxonomy alignment between predicted basic emotions and epistemic labels, (2) dimensional valence–arousal structure, and (3) cross-model agreement under domain shift.

4.1. Cross-Taxonomy Alignment

Figure 3 shows the row-normalized confusion matrix between epistemic labels and OpenFace 3.0 predictions. We see a consistent pattern across all epistemic states: *Neutral* is the most frequent prediction. For example, 43.8% of *Curious* instances, 33.3% of *Confused* instances, and 37.5% of *Frustrated* instances are predicted as *Neutral*. No epistemic label shows a strong one-to-one mapping to any basic emotion category. This pattern persists across all evaluated categorical models. Table 2 summarizes, for each epistemic label, the most common predicted basic emotion and its corresponding proportion. Across models, *Neutral* is overwhelmingly the dominant prediction for nearly all epistemic states. In the few cases where another label appears most frequently (e.g., *Happy* or *Sad*), the mapping is inconsistent across models and lacks an interpretable semantic relationship to the epistemic construct.

4.2. Valence–Arousal Structure

We next analyze continuous valence–arousal (V–A) outputs. Figure 4 visualizes the distribution of epistemic labels in V–A space. We observe that epistemic states do not form clearly separable regions, as the scatter plot reveal substantial overlap across labels, showing large within-label variance relative to between-label differences. When grouped by predicted *basic emotion* categories, we observe

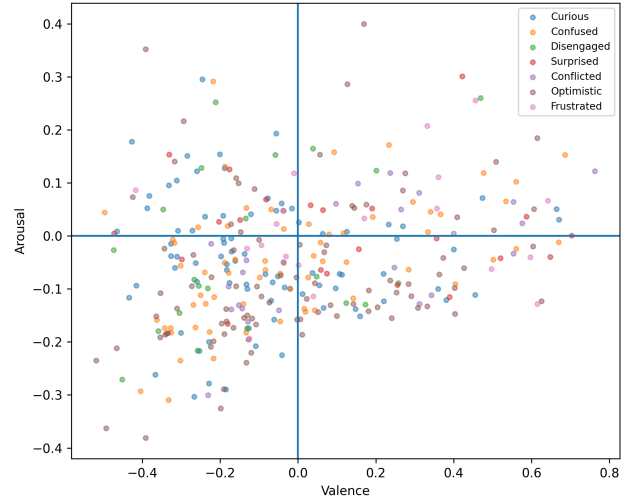


Figure 4. Valence–arousal scatter plot grouped by epistemic labels. Substantial overlap and high within-label variance are observed, with no clearly separable clusters.

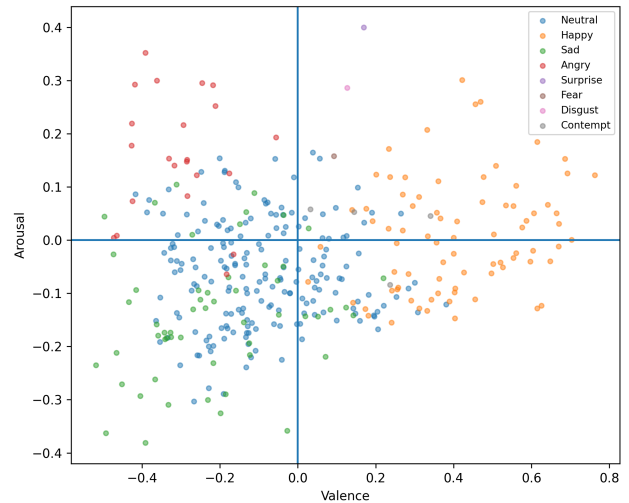


Figure 5. Valence–arousal scatter plot grouped by predicted basic emotion categories. Distinct clusters emerge (e.g., *Happy* in high valence, *Angry* in low valence/high arousal), indicating internal coherence of the dimensional model with respect to its training taxonomy.

a more clear structure (Figure 5). For example, *Happy* predictions cluster in high-valence regions, *Sad* in low-valence and lower-arousal regions, and *Angry* in low-valence but higher-arousal regions.

4.3. Cross-Model Agreement

Finally, we assess robustness across architectures under domain shift by analyzing pairwise confusion matrices between categorical predictions generated by OpenFace 3.0 and LibreFace as representative models on subsets of the

Table 2. Most frequent predicted basic emotion for each epistemic label and their proportion across categorical models.

Epistemic Label	OpenFace 3.0	LibreFace	POSTER++	DDAMFN	EmotiEffLib	Qwen 2.5-VL
Curious	Neutral (0.44)	Neutral (0.51)	Neutral (0.73)	Neutral (0.73)	Neutral (0.70)	Neutral (0.96)
Confused	Neutral (0.33)	Neutral (0.43)	Neutral (0.54)	Neutral (0.54)	Neutral (0.44)	Neutral (0.93)
Disengaged	Neutral (0.35)	Neutral (0.59)	Neutral (0.71)	Neutral (0.71)	Neutral (0.48)	Neutral (0.81)
Surprised	Neutral (0.31)	Neutral (0.53)	Neutral (0.38)	Happy (0.50)	Neutral (0.41)	Neutral (0.95)
Conflicted	Neutral (0.39)	Happy (0.52)	Neutral (0.43)	Neutral (0.43)	Neutral (0.46)	Neutral (0.88)
Optimistic	Neutral (0.37)	Sad (0.46)	Neutral (0.60)	Neutral (0.60)	Neutral (0.55)	Neutral (0.97)
Frustrated	Neutral (0.38)	Neutral (0.44)	Neutral (0.56)	Neutral (0.56)	Neutral (0.47)	Neutral (0.84)

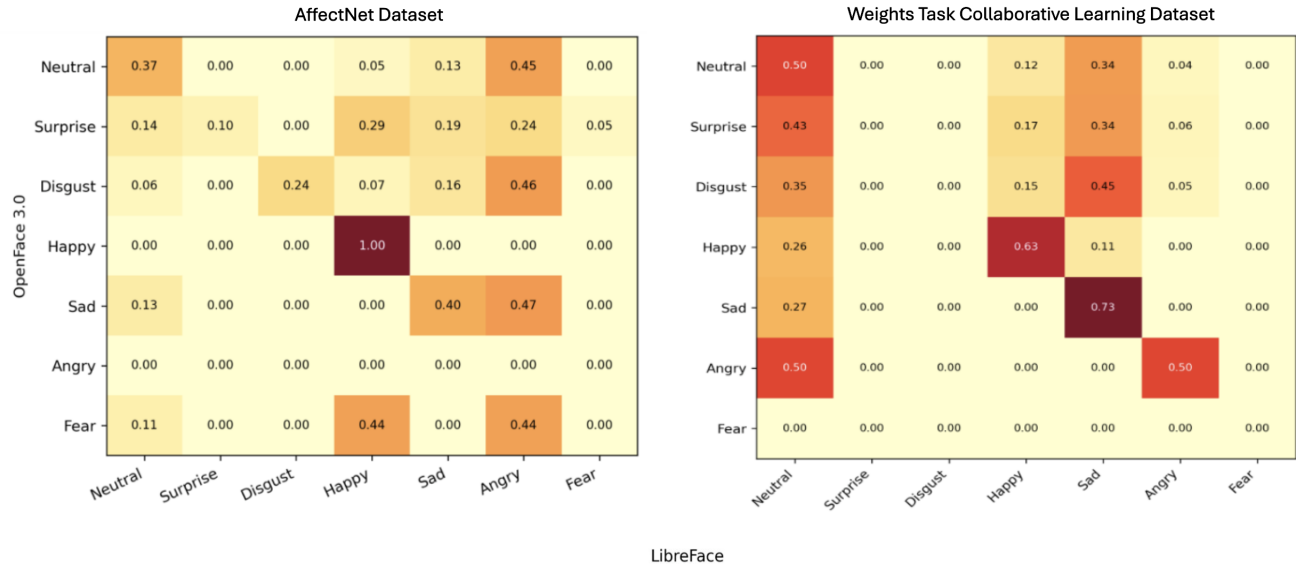


Figure 6. Cross-model confusion matrices between OpenFace 3.0 and LibreFace predictions on AffectNet (left) and EECPS-WT (right).

454 AffectNet and EECPS-WT dataset. Figure 6 presents the
455 two confusion matrices.

456 For AffectNet, agreement between models varies sub-
457 stantially across emotion categories. The strongest align-
458 ment occurs for *Happy*, where predictions coincide almost
459 perfectly between the two models. Moderate agreement is
460 also observed for *Sad*. However, other categories exhibit
461 substantial disagreement. For example, images predicted as
462 *Neutral* by OpenFace 3.0 are frequently labeled as *Angry*
463 or *Sad* by LibreFace, and *Disgust* predictions often corre-
464 spond to *Angry*. Despite these inconsistencies, the Affect-
465 Net matrix exhibits a relatively diverse distribution of pre-
466 dicted labels, including *Fear*, *Disgust*, and *Angry*, reflecting
467 the presence of more prototypical facial expressions within
468 the benchmark dataset.

469 In contrast to AffectNet, predictions in EECPS-WT are
470 heavily concentrated in a small subset of categories, partic-
471 ularly *Neutral*, *Sad*, and *Happy*. Other categories such as
472 *Fear*, *Disgust*, and *Angry* appear rarely or not at all. Addi-

tionally, cross-model agreement is limited, with predictions
frequently distributed across multiple categories rather than
aligning along the diagonal.

5. Discussion

Our findings highlight three interrelated challenges in de-
ploying pretrained FER systems in collaborative educa-
tional settings.

5.1. Basic Emotions and Epistemic States

Across models, epistemic states are predominantly mapped
to *Neutral*. This collapse suggests that models trained on
web-scale benchmarks appear calibrated toward canonical
and often exaggerated facial prototypes and default to *Neu-
tral* when subtle or ambiguous expressions are encountered.
In collaborative problem-solving (CPS) settings, where af-
fect may be regulated, masked, or embedded within task-
focused expressions, such calibration leads to a collapse of
predictions into a small subset of categories. Such emo-

490 tions may instead manifest as mild brow furrows, brief gaze
491 shifts, or micro-level muscular changes that fall below the
492 decision thresholds learned from prototypical training data.

493 5.2. Valence–Arousal and Epistemic States

494 Dimensional valence–arousal modeling provides a richer
495 representational space, yet epistemic states remain highly
496 overlapping in V–A space. While the dimensional model
497 behaves consistently with respect to canonical basic emo-
498 tions, it does not yield separable structure for collabora-
499 tively reported epistemic labels. This contrast indicates that
500 the dimensional model is internally coherent with respect
501 to the basic emotion taxonomy on which it was trained.
502 However, that structure does not transfer to epistemic states
503 observed in CPS. Dimensional representations therefore in-
504 herit the same domain alignment limitations as categorical
505 models.

506 5.3. Cross-Model Reliability

507 The cross-model agreement analysis reveals several im-
508 portant insights about the reliability of current FER sys-
509 tems when applied to educational contexts. First, even
510 when evaluated on AffectNet images—data drawn from
511 a large-scale facial expression benchmark—agreement be-
512 tween OpenFace 3.0 and LibreFace is inconsistent across
513 several emotion categories. While both models reliably
514 identify highly prototypical expressions such as *Happy*, pre-
515 dictions diverge for other emotions including *Disgust*, *An-*
516 *gry*, and *Neutral*. This suggests that different model archi-
517 tectures may learn distinct feature representations and de-
518 cision boundaries despite sharing a common emotion tax-
519 onomy. Second, cross-model disagreement becomes even
520 more pronounced when models are applied to CPS. In
521 EECPS-WT, predictions from both models collapse into
522 a small subset of categories—primarily *Neutral*, *Sad*, and
523 *Happy*. Emotions such as *Fear*, *Disgust*, and *Angry*, which
524 appear in AffectNet predictions, are rarely produced in the
525 collaborative setting. This pattern suggests that facial ex-
526 pressions observed during CPS may not strongly resem-
527 ble the canonical facial configurations present in benchmark
528 datasets.

529 These results are particularly noteworthy given that both
530 OpenFace 3.0 and LibreFace were trained using datasets de-
531 rived in part from AffectNet. However, the two systems
532 differ substantially in their training pipelines and additional
533 training data. LibreFace also incorporates pretraining on
534 the FFHQ [37] and EmotioNet [4] datasets and fine-tuning
535 on DISFA [38], while OpenFace 3.0 employs a multi-task
536 learning architecture trained to jointly predict facial land-
537 marks, action units, and emotion logits. Differences in
538 training objectives and auxiliary datasets may therefore in-
539 fluence how each model interprets facial configurations in
540 naturalistic environments.

541 For educational applications, these findings raise impor-
542 tant concerns about measurement reliability. If two widely
543 used FER systems produce different predictions when ap-
544 plied to the same student behavior, it becomes difficult to
545 interpret model outputs as stable indicators of learner affect.
546 In CPS settings where emotional signals may be subtle, reg-
547 ulated, or embedded within task-focused activity, reliance
548 on facial expressions alone may therefore provide an in-
549 complete or inconsistent representation of learners’ internal
550 states. This highlights the need for multimodal approaches
551 and domain-specific evaluation protocols when deploying
552 affective computing systems in educational research.

553 6. Limitations

554 This study has several limitations. First, our analysis
555 is based on a single collaborative problem-solving (CPS)
556 dataset collected within one institutional context with a rel-
557 atively small participant pool. Further validation across di-
558 verse classrooms, age groups, and cultural contexts is nec-
559 essary to assess the generalizability of our findings. Sec-
560 ond, we evaluate a limited set of FER models that primarily
561 operate on static images. Future work should extend this
562 analysis to multimodal and temporally-aware architectures,
563 including video-based models and vision–language models
564 fine-tuned for affect recognition [39]. Finally, our dataset
565 relies on retrospective cued-recall annotations, which may
566 be susceptible to recall bias. Future work should explore
567 complementary annotation strategies, including hybrid ap-
568 proaches that integrate behavioral and physiological sig-
569 nals.

570 7. Conclusion

571 We evaluated web-trained FER systems in a collabora-
572 tive problem-solving (CPS) setting annotated with epis-
573 temic states via retrospective cued-recall. Across multi-
574 ple architectures, categorical models largely mapped epis-
575 temic states to *Neutral*, while valence–arousal representa-
576 tions failed to produce separable structure among labels.
577 Low inter-model agreement further indicated instability un-
578 der domain shift. Dimensional predictions remained inter-
579 nally consistent with canonical basic emotions, suggesting
580 the misalignment arises not from model failure but from
581 a mismatch between web-trained emotion taxonomies and
582 collaboratively expressed epistemic affect. These results
583 suggest that these modern FER systems trained on web-
584 scale datasets may not directly transfer to CPS contexts, and
585 their outputs should be interpreted cautiously when used in
586 human-centered AI research.

587 References

- 588 [1] A. Mollahosseini, B. Hasani, and M. H. Mahoor, “Affect-
589 net: A database for facial expression, valence, and arousal

- 590 computing in the wild,” *IEEE transactions on affective computing*, vol. 10, no. 1, pp. 18–31, 2017. 1, 2, 3, 5
- 591
- 592 [2] T. Baltrušaitis, P. Robinson, and L.-P. Morency, “Openface:
593 an open source facial behavior analysis toolkit,” in *2016
594 IEEE winter conference on applications of computer vision
595 (WACV)*. IEEE, 2016, pp. 1–10. 2
- 596 [3] J. Hu, L. Mathur, P. P. Liang, and L.-P. Morency, “Openface
597 3.0: A lightweight multitask system for comprehensive facial
598 behavior analysis,” in *2025 IEEE 19th International Confer-
599 ence on Automatic Face and Gesture Recognition (FG)*.
600 IEEE, 2025, pp. 1–11. 2, 4
- 601 [4] C. Fabian Benitez-Quiroz, R. Srinivasan, and A. M. Mar-
602 tinez, “Emotionet: An accurate, real-time algorithm for the
603 automatic annotation of a million facial expressions in the
604 wild,” in *Proceedings of the IEEE conference on computer
605 vision and pattern recognition*, 2016, pp. 5562–5570. 2, 3, 8
- 606 [5] J. Mao, R. Xu, X. Yin, Y. Chang, B. Nie, A. Huang, and
607 Y. Wang, “Poster++: A simpler and stronger facial expres-
608 sion recognition network,” *Pattern Recognition*, vol. 157, p.
609 110951, 2025. 2, 4
- 610 [6] K. Ahuja, D. Kim, F. Xhakaj, V. Varga, A. Xie, S. Zhang,
611 J. E. Townsend, C. Harrison, A. Ogan, and Y. Agarwal,
612 “Edusense: Practical classroom sensing at scale,” *Proceed-
613 ings of the ACM on Interactive, Mobile, Wearable and Ubiqu-
614 itous Technologies*, vol. 3, no. 3, pp. 1–26, 2019. 2
- 615 [7] P. Goldberg, Ö. Sümer, K. Stürmer, W. Wagner, R. Göllner,
616 P. Gerjets, E. Kasneci, and U. Trautwein, “Attentive or not?
617 toward a machine learning approach to assessing students’
618 visible engagement in classroom instruction,” *Educational
619 Psychology Review*, vol. 33, no. 1, pp. 27–49, 2021. 3
- 620 [8] A. Joshi, D. Alessio, J. Magee, J. Whitehill, I. Arroyo,
621 B. Woolf, S. Sclaroff, and M. Betke, “Affect-driven learn-
622 ing outcomes prediction in intelligent tutoring systems,” in
623 *2019 14th IEEE international conference on automatic face
624 & gesture recognition (fg 2019)*. IEEE, 2019, pp. 1–5.
- 625 [9] X. Tang, Y. Gong, Y. Xiao, J. Xiong, and L. Bao, “Facial
626 expression recognition for probing students’ emotional en-
627 gagement in science learning,” *Journal of Science Education
628 and Technology*, vol. 34, no. 1, pp. 13–30, 2025. 2
- 629 [10] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville,
630 M. Mirza, B. Hamner, W. Cukierski, Y. Tang, D. Thaler,
631 D.-H. Lee *et al.*, “Challenges in representation learning: A
632 report on three machine learning contests,” in *International
633 conference on neural information processing*. Springer,
634 2013, pp. 117–124. 2, 3
- 635 [11] A. P. Fard, M. M. Hosseini, T. D. Sweeny, and M. H. Ma-
636 hoor, “Affectnet+: A database for enhancing facial expres-
637 sion recognition with soft-labels,” *IEEE Transactions on Af-
638 fective Computing*, 2025. 2, 4
- 639 [12] J. A. Russell, “A circumplex model of affect,” *Journal of per-
640 sonality and social psychology*, vol. 39, no. 6, p. 1161, 1980.
641 2, 3
- 642 [13] F. Zhou, S. Kong, C. C. Fowlkes, T. Chen, and B. Lei, “Fine-
643 grained facial expression analysis using dimensional emo-
644 tion model,” *Neurocomputing*, vol. 392, pp. 38–49, 2020. 2
- 645 [14] S. D’Mello, R. W. Picard, and A. Graesser, “Toward an
646 affect-sensitive autotutor,” *IEEE Intelligent Systems*, vol. 22,
647 no. 4, pp. 53–61, 2007. 2, 3
- [15] R. S. Baker, S. K. D’Mello, M. M. T. Rodrigo, and A. C. 648
Graesser, “Better to be frustrated than bored: The inci- 649
dence, persistence, and impact of learners’ cognitive– 650
affective states during interactions with three different 651
computer-based learning environments,” *International Jour- 652
nal of Human-Computer Studies*, vol. 68, no. 4, pp. 223–241, 653
2010. 3 654
- [16] E. B. Cloude, A. Munshi, J. A. Andres, J. Ocumpaugh, R. S. 655
Baker, and G. Biswas, “Exploring confusion and frustra- 656
tion as non-linear dynamical systems,” in *Proceedings of the 657
14th learning analytics and knowledge conference*, 2024, pp. 658
241–252. 2 659
- [17] R. A. Calvo and S. D’Mello, “Affect detection: An interdis- 660
ciplinary review of models, methods, and their applications,” 661
IEEE Transactions on affective computing, vol. 1, no. 1, pp. 662
18–37, 2010. 2 663
- [18] D. M. Russell and M. Oren, “Retrospective cued recall: a 664
method for accurately recalling previous user behaviors,” in 665
*2009 42nd Hawaii International Conference on System Sci- 666
ences*. IEEE, 2009, pp. 1–9. 2, 3 667
- [19] S. Anindho, V. Venkatesha, and N. Blanchard, “A method- 668
ological framework for capturing cognitive-affective states 669
in collaborative learning,” *arXiv preprint arXiv:2507.01166*, 670
2025. 2, 3 671
- [20] S. Li, W. Deng, and J. Du, “Reliable crowdsourcing and deep 672
locality-preserving learning for expression recognition in the 673
wild,” in *Proceedings of the IEEE conference on computer 674
vision and pattern recognition*, 2017, pp. 2852–2861. 2, 3 675
- [21] S. Zhang, Y. Zhang, Y. Zhang, Y. Wang, and Z. Song, “A 676
dual-direction attention mixed feature network for facial ex- 677
pression recognition,” *Electronics*, vol. 12, no. 17, p. 3595, 678
2023. 2, 4 679
- [22] D. Chang, Y. Yin, Z. Li, M. Tran, and M. Soleymani, “Li- 680
breface: An open-source toolkit for deep facial expression 681
analysis,” in *Proceedings of the IEEE/CVF winter confer- 682
ence on applications of computer vision*, 2024, pp. 8205– 683
8215. 2, 4 684
- [23] G. A. Miller, “Wordnet: a lexical database for english,” *Com- 685
munications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995. 2 686
- [24] M. K. Underwood, “Peer social status and children’s under- 687
standing of the expression and control of positive and nega- 688
tive emotions,” *Merrill-Palmer Quarterly (1982-)*, pp. 610– 689
634, 1997. 2 690
- [25] L. Isenbarger and M. Zembylas, “The emotional labour of 691
caring in teaching,” *Teaching and teacher education*, vol. 22,
no. 1, pp. 120–134, 2006. 2 692
693
- [26] P. Ekman, “An argument for basic emotions,” *Cognition & 694
emotion*, vol. 6, no. 3–4, pp. 169–200, 1992. 3 695
- [27] M. P. Cross, A. M. Acevedo, and J. F. Hunter, “A critique of 696
automated approaches to code facial expressions: What do 697
researchers need to know?” *Affective Science*, vol. 4, no. 3,
pp. 500–505, 2023. 3 698
699
- [28] A. Horvers, N. Tombeng, T. Bosse, A. W. Lazonder, and 700
I. Molenaar, “Detecting emotions through electrodermal ac- 701
tivity in learning contexts: A systematic review,” *Sensors*,
vol. 21, no. 23, p. 7869, 2021. 3 702
703

- 704 [29] J. W. Y. Chung, H. C. F. So, M. M. T. Choi, V. C. M. Yan,
705 and T. K. S. Wong, "Artificial intelligence in education: Us-
706 ing heart rate variability (hrv) as a biomarker to assess emo-
707 tions objectively," *Computers and Education: Artificial In-*
708 *telligence*, vol. 2, p. 100011, 2021. 3
- 709 [30] I. Khebour, R. Brutti, I. Dey, R. Dickler, K. Sikes, K. Lai,
710 M. Bradford, B. Cates, P. Hansen, C. Jung *et al.*, "When text
711 and speech are not enough: A multimodal dataset of collab-
712 oration in a situated task," *Journal of open humanities data*,
713 vol. 10, 2024. 3
- 714 [31] S. Anindho, V. Venkatesha, M. Bradford, A. M. Cleary, and
715 N. Blanchard, "An exploration of internal states in collab-
716 orative problem solving," in *International Conference on*
717 *Human-Computer Interaction*. Springer, 2025, pp. 135–
718 150. 4
- 719 [32] J. Deng, J. Guo, E. Ververas, I. Kotsia, and S. Zafeiriou,
720 "Retinaface: Single-shot multi-level face localisation in the
721 wild," in *Proceedings of the IEEE/CVF conference on com-*
722 *puter vision and pattern recognition*, 2020, pp. 5203–5212.
723 4
- 724 [33] A. Savchenko, "Facial expression recognition with adaptive
725 frame rate based on multiple testing correction," in *Pro-*
726 *ceedings of the 40th International Conference on Machine*
727 *Learning (ICML)*, ser. Proceedings of Machine Learning
728 Research, A. Krause, E. Brunskill, K. Cho, B. Engel-
729 hardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR,
730 23–29 Jul 2023, pp. 30 119–30 129. [Online]. Available:
731 <https://proceedings.mlr.press/v202/savchenko23a.html> 4, 5
- 732 [34] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen,
733 X. Liu, J. Wang, W. Ge *et al.*, "Qwen2-vl: Enhancing vision-
734 language model's perception of the world at any resolution,"
735 *arXiv preprint arXiv:2409.12191*, 2024. 4
- 736 [35] S. I. Serengil and A. Ozpinar, "Hyperextended lightface: A
737 facial attribute analysis framework," in *2021 International*
738 *Conference on Engineering and Emerging Technologies*
739 *(ICEET)*. IEEE, 2021, pp. 1–4. [Online]. Available:
740 <https://ieeexplore.ieee.org/document/9659697> 5
- 741 [36] D. Kollias and S. Zafeiriou, "Aff-wild2: Extending the
742 aff-wild database for affect recognition," *arXiv preprint*
743 *arXiv:1811.07770*, 2018. 5
- 744 [37] T. Karras, S. Laine, and T. Aila, "A style-based generator ar-
745 chitecture for generative adversarial networks," in *Proceed-*
746 *ings of the IEEE/CVF conference on computer vision and*
747 *pattern recognition*, 2019, pp. 4401–4410. 8
- 748 [38] S. M. Mavadati, M. H. Mahoor, K. Bartlett, P. Trinh, and
749 J. F. Cohn, "Disfa: A spontaneous facial action inten-
750 sity database," *IEEE Transactions on Affective Computing*,
751 vol. 4, no. 2, pp. 151–160, 2013. 8
- 752 [39] A. Chaubey, X. Guan, and M. Soleymani, "Face-llava: Facial
753 expression and attribute understanding through instruction
754 tuning," in *Proceedings of the IEEE/CVF Winter Conference*
755 *on Applications of Computer Vision*, 2026, pp. 2648–2660.
756 8