

# CONTEXTUALLY-RICH HUMAN AFFECT PERCEPTION USING MULTIMODAL SCENE INFORMATION

Digbalay Bose, Rajat Hebbar, Krishna Somandepalli <sup>†</sup>, Shrikanth Narayanan

Signal Analysis and Interpretation Laboratory, University of Southern California, CA 90089

dbose@usc.edu, rajathebb@usc.edu, somandep@usc.edu, shri@ee.usc.edu

## ABSTRACT

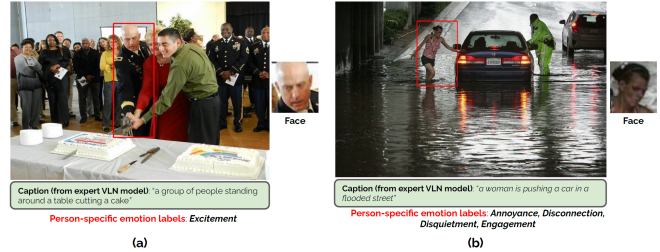
The process of human affect understanding involves the ability to infer person specific emotional states from various sources including images, speech, and language. Affect perception from images has predominantly focused on expressions extracted from salient face crops. However, emotions perceived by humans rely on multiple contextual cues including social settings, foreground interactions, and ambient visual scenes. In this work, we leverage pretrained vision-language (VLN) models to extract descriptions of foreground context from images. Further, we propose a multimodal context fusion (MCF) module to combine foreground cues with the visual scene and person-based contextual information for emotion prediction. We show the effectiveness of our proposed modular design on two datasets associated with natural scenes and TV shows.

**Index Terms**— Emotion recognition, Context understanding, Multimedia, Multimodal vision-language pretrained models, Multimodal interaction modeling

## 1. INTRODUCTION

There has been increased interest in understanding the affective processes associated with various facets of human emotions [1]. An integral part of affective understanding is the ability to infer expressions of human emotions from various sources like images [2], speech [3] and language use. Affect recognition systems have enabled multiple human-centered applications, notably in healthcare (depression detection [4], autism-spectrum diagnosis [5]) and learning [6].

Affect recognition from images has largely focused on facial expressions [7, 8] along a fixed set of categories. Moreover, facial expression based methods typically consider crops of a single face, which might provide ambiguous signals for classifying perceived emotions. Emotion perception in humans typically relies on multimodal behavioral cues, that go beyond facial expressions, such as voice and language [9]. However, are there additional contextual cues beyond behavioral expressions, such as of face and language, that mediate human emotion perception? Studies have shown that contextual information including the social setting, interaction type, and ambient location can play a key role [10]. Context in images is driven by visual scenes [11] or specific locations such as outdoor, indoor, kitchen, living room etc and the interactions between the various entities in the scene. An example is shown in Fig.1 with respect to both (a) positive and (b) negative emotions. In Fig.1 (a), the face crop provides a negative signal whereas the overall scene including the generated caption from a pretrained vision-language model OFA [12] indicates a positive event associated with the person. In



**Fig. 1.** Examples from EMOTIC dataset showing the importance of context for estimating person specific emotion labels. Extracted captions from a pretrained VLN model (OFA) capture the foreground contexts.

Fig. 1(b), while the face crop provides a noisy, incomplete signal for perceiving the expressed emotional state, the overall context of the visual scene plus its descriptive caption indicates the distressing situation associated with street flooding.

Recent advances in multimodal vision-language (VLN) pretraining [13, 12] has resulted in task-agnostic transformer-based models that can be used for variety of tasks such as image captioning and visual question answering. In this work, we employ the pretrained VLN models as experts for describing foreground context in terms of captions. Further we consider contextual information in terms of the individual persons (whole-body appearance or face) and the visual scenes. In order to effectively leverage multiple contextual sources we propose attention-based multimodal context fusion (MCF) module to predict both discrete emotion labels and continuous-valued arousal, valence and dominance values. We show the effectiveness of our methods on two publicly available datasets: EMOTIC [14] (natural scenes) and CAER-S [15](TV shows).

## 2. RELATED WORK

**Context in Emotion Perception:** The role of context extraneous to a person (beyond their traits and behavior) in the perception of their expressed emotion has been studied from the perspective of scenes [10], and cultures [16]. In [17], the perceivable-encoding context and the prior knowledge available with the perceivers are reported as the major sources of context for influencing emotion perception. Situational context like reactions from other people has been considered in [18] as a means to decode emotional states of persons in consideration.

**Context-based image datasets:** Image-based emotion recognition datasets like AffectNet [19], FER [20], DFEW [7] primarily focus

<sup>†</sup>The work was done while the author was at USC

on signals encoded in facial expressions. Since emotion perception depends on where the facial configuration is present, datasets like EMOTIC [14] and CAER [15] have been proposed to incorporate contextual information in terms of visual scenes and social interactions. In the case of EMOTIC, annotators have marked person instances in unconstrained environments with apparent emotional states based on the existing scene context. However, the annotation process in CAER revolves around TV shows with primary focus on interactions-driven context.

**Context modeling approaches:** [14, 15] explore context modeling in terms of dual stream models for processing body and whole image streams. [21] also uses a dual stream network with context stream, modeled using an affective graph composed of region proposals. [22] uses depth and pose as additional contextual signals with existing streams like scene, face for predicting person specific emotion. [23] explores contextual modeling in short movie clips by considering scene and action characteristics along with body (including face) based signals. In contrast, our approach uses natural language descriptions to describe the foreground context along with scene and person specific streams.

**Multimodal vision-language models:** Vision language (VLN) models like OFA [12], VL-T5 [13], ALBEF [24] are pretrained on large-scale image-text pairs curated from the web, thus enabling its usage in diverse tasks like image captioning, retrieval, visual-question answering etc. In our formulation, we harness the capabilities of VLN models to generate descriptive captions since they contain condensed description of the foreground context in terms of entities including persons.

### 3. PROBLEM FORMULATION

Given an image  $I$  and a person-specific context in the form of bounding box  $[x, y, w, h]$ , the task is to predict the emotional state  $p$  associated with a person as  $p_{disc}, p_{cont} = F(I, [x, y, w, h])$ .  $p_{disc}$  and  $p_{cont}$  refer to the predicted set of discrete emotion categories and continuous arousal, valence and dominance values, respectively.  $F$  refers to the deep neural network used for estimating the discrete and continuous affective states. The design of  $F$  is dependent on extraction of multiple contextual information from the given image  $I$ , that are listed as below:

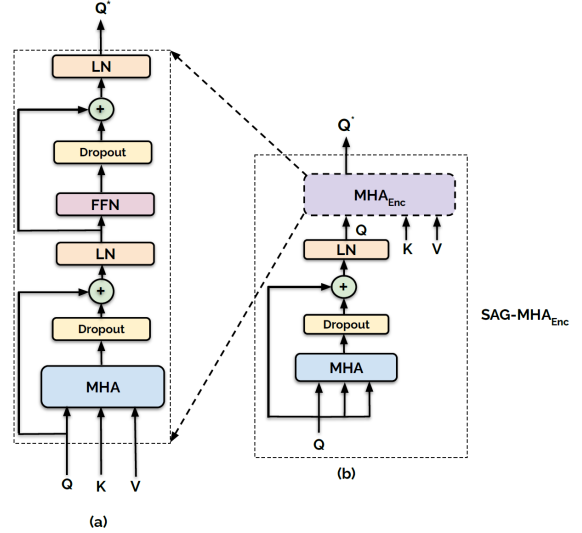
**Visual scene context:** The underlying visual scene ( $VS$ ) (e.g., kitchen, bar, football field etc) plays a role in influencing the emotional state of a person. Here we use a ViT [25] model ( $f_{VS}$ ) finetuned on Places365 [26] as the backbone network for extracting visual scene representations ( $e_{VS}$ ) from  $I$ .

**Person context:** The person-specific context is extracted using a whole-body or facial bounding box, denoted by  $[x, y, w, h]$  from image  $I$ . The cropped person instance is passed through a person encoder ( $f_{PE}$ ) i.e. Resnet34 [27] for extracting person-centric representations ( $e_{PE}$ ).

**Language driven foreground context:** Natural language description of image  $I$  provides foreground (FG) context in terms of entities including persons and their interactions. We use a 12-layer transformer encoder-decoder model  $OFA_{large}$  [12] as  $f_{expert}$  to extract the foreground specific captions for image  $I$ . For extracting text representations ( $e_{FG}$ ) of the captions, we use BERT's [28] pretrained encoder ( $f_{FG}$ ) from HuggingFace [29].

### 4. MULTIMODAL CONTEXT FUSION (MCF) MODULE

The multimodal context fusion module is composed of two parallel streams, associated with foreground and visual scene based contexts.



**Fig. 2.** Outline of (a)  $MHA_{enc}$  (b)  $SAG-MHA_{enc}$  layers. **LN:** Layer Norm, **FFN:** Feedforward network, **MHA:** Multihead attention, **SAG:** Self Attention Augmented

The basic operation in individual streams is a cross-modal encoder block  $CM_{enc}$  composed of  $L$  encoder layers. As shown in Fig 2, we consider two designs for the encoder layer i.e.,  $MHA_{enc}$  and  $SAG-MHA_{enc}$ . The set of operations in encoder  $MHA_{enc}$  layer for query ( $Q$ ), key ( $K$ ) and value ( $V$ ) representations are listed as follows:

$$\begin{aligned} Q' &= LN(Q + Dropout(MHA(Q, K, V))) \\ Q^* &= LN(Dropout(FFN(Q')) + Q') \end{aligned} \quad (1)$$

Here  $MHA$ ,  $LN$  and  $FFN$  refer to Multi-head attention, layer-norm operation and feed-forward neural network respectively. The  $SAG-MHA_{enc}$  layer consists of a multi-head attention based transformation of the query representations followed by input to the  $MHA_{enc}$  layer. The design of  $SAG-MHA_{enc}$  is inspired from multimodal co-attention layer proposed in [30] for visual question answering task.

$$\begin{aligned} Q' &= LN(Q + Dropout(MHA(Q, Q, Q))) \\ Q^* &= MHA_{enc}(Q', K, V) \end{aligned} \quad (2)$$

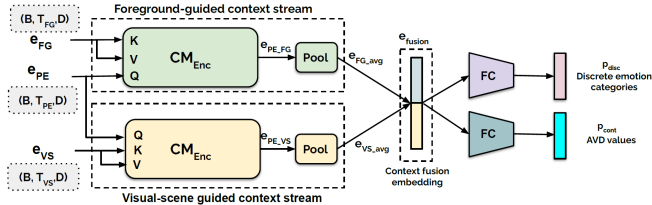
In  $CM_{enc}$ , the output from the  $i$ th encoder layer  $Enc_i$  is passed as query ( $Q$ ) to the subsequent layer with the key ( $K$ ) and value ( $V$ ) remaining the same. Here,  $Enc_i$  layer can be either  $MHA_{enc}$  or  $SAG-MHA_{enc}$ .

$$\begin{aligned} Q_i &= Enc_i(Q_{i-1}, K, V) \quad i > 0 \\ Q_0 &= Enc_0(Q, K, V) \end{aligned} \quad (3)$$

We use separate  $CM_{enc}$  blocks for processing the foreground and visual scene guided context streams. MCF ( $MHA_{enc}$ ) and MCF ( $SAG-MHA_{enc}$ ) consists of 4  $MHA_{enc}$  layers (8 heads and hidden dimension=512) and 3  $SAG-MHA_{enc}$  layers (8 heads and hidden dimension=768) in the  $CM_{enc}$  blocks respectively. The details of the respective context-guided streams are listed as follows:

**Foreground-guided context stream:** We use  $e_{PE}$  as query ( $Q$ ) and  $e_{FG}$  as key and value inputs to the  $CM_{enc}$  encoder.

**Visual-scene guided context stream:** Similar to the foreground-guided context stream, we use  $e_{PE}$  as query ( $Q$ ) and  $e_{VS}$  as key



**Fig. 3.** Outline of the *MCF* module. Separate  $CM_{enc}$  blocks are used to model the foreground and visual scene context streams.  $B$ : Batch size.  $D$ : Hidden dimension. Token length of  $T_{FG}$ : text representations,  $T_{PE}$ : person representations,  $T_{VS}$ : visual scene representations. AVD: Arousal, Valence, Dominance

and value inputs to the  $CM_{enc}$  encoder.

**Context fusion:** The outputs from the context-guided streams i.e.,  $e_{PE,FG}$  and  $e_{PE,VS}$  are average pooled and concatenated to obtain a fused embedding as:  $e_{fusion} = [e_{FG,avg}; e_{VS,avg}]$

## 5. EXPERIMENTS

**EMOTIC:** We use a non-intersecting split of 13584 and 4389 images for training and validation. For testing we use the publicly available split of 5108 images. For joint prediction of 26 discrete emotion classes and the continuous-valued AVD ratings, we use multiple fully-connected (FC) heads in the *MCF* module with  $e_{fusion}$  as input (Fig 3). The person specific instance in each image is defined by the ground truth person box. We do not consider face as a part of person-specific context since approx 25% of images do not have visible faces.

For training *MCF* with  $MHA_{enc}$  and  $SAG-MHA_{enc}$  layers, we use AdamW [31] ( $lr=2e-5$ ) and Adam [32] ( $lr=2e-4$ ,  $exp(\gamma=0.90)$ ) with batch sizes 32 and 64 respectively.<sup>1</sup> We use SGD ( $lr=1e-2$ ,  $exp(\gamma=0.90)$ ) with a batch size of 64 while training the person-crop only Resnet34 model ( $PO_{R34}$ ) in Table 3. For training all the models associated with EMOTIC, we use a weighted combination of binary-cross entropy (BCE) and mean squared error (MSE) losses.

$$Loss = \lambda_1 BCE(p_{disc}, y_{disc}) + \lambda_2 MSE(p_{cont}, y_{cont}) \quad (4)$$

Here  $y_{disc}$  and  $y_{cont}$  refer to ground truth discrete emotion labels and continuous arousal valence dominance ratings. The optimal weights  $\lambda_1$  and  $\lambda_2$  are tuned using the validation split.

**CAER-S:** We use a non-intersecting split of 39099 and 9769 video frames across 79 TV shows for training and validation. For testing we use the public split of 20913 video frames. Since face is a dominant signal for persons in TV shows, we use MTCNN [33]<sup>2</sup> to obtain face crops. We have a single fully-connected (FC) head with  $e_{fusion}$  as input for predicting 7 discrete emotion classes.

For training *MCF* with  $MHA_{enc}$  and  $SAG-MHA_{enc}$  layers, we use Adam ( $lr=2e-4$ ,  $exp(\gamma=0.90)$ ) with a batch size of 64. We use Adam ( $lr=1e-4$ ,  $exp(\gamma=0.75)$ ) with a batch size of 64 while training the face-crop only Resnet34 model ( $FO_{R34}$ ) in Table 4. For training all the models associated with CAER-S, we use multi-class cross entropy loss.

We conduct our experiments using the Pytorch [34] library. We set maximum sequence length  $T_{FG}$  as 512 for the captions (Fig 3).

<sup>1</sup>exp is exponential scheduler

<sup>2</sup><https://github.com/timesler/facenet-pytorch>

For visual scene and person representations we use  $T_{VS} = 197$  and  $T_{PE} = 49$  respectively.

## 6. RESULTS

### 6.1. Comparison with state of the art

We compare performance of *MCF* (Enc) under two settings where Enc refers to the encoder layer used i.e.,  $MHA_{enc}$  and  $SAG-MHA_{enc}$  with existing methods in Table 1. We can see that *MCF* under both settings performs better than prior methods like [14], [21], and [15] that rely on dual stream (person + whole image approach) and do not use explicit pose information. Furthermore, in contrast to previous methods, we consider language driven foreground (captions) and visual scene contexts instead of end-to-end modeling of whole image based information. For a fair comparison with [22], the current *MCF* framework can be potentially expanded to include other person specific streams like face and explicit pose information. From Table 2, we can see that in CAER-S, a fully finetuned

Model	mAP
Kosti et. al [14]	27.38
Zhang et. al [21]	28.42
Lee et.al [15]	20.84
<b>MCF (<math>MHA_{enc}</math>)</b>	<b>29.53 (0.001)</b>
MCF ( $SAG-MHA_{enc}$ )	28.58 (0.003)
EmotiCon [22]	32.03

**Table 1.** Comparison of *MCF* module with state of the art in EMOTIC (Test). **mAP:** mean average precision. Average of runs with 5 random seeds reported with standard deviation for *MCF*.

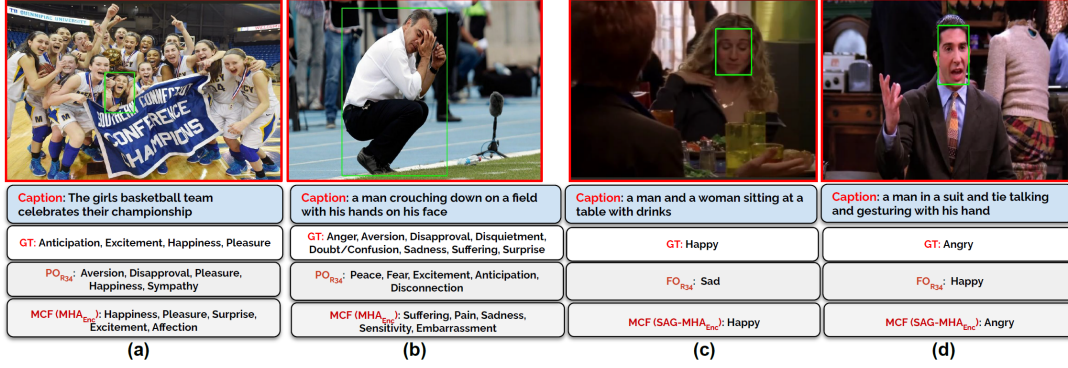
Resnet34 model trained on face crops ( $FO_{R34}$ ) obtains a high accuracy of 77.35 since facial expressions provide dominant signals for emotion classification in TV shows. However, inclusion of both foreground context through captions and visual scene information in *MCF* ( $MHA_{enc}$ ) results in better performance (**79.63**) as compared to both  $FO_{R34}$  and baseline attention fusion method CAER-Net-S.

Model	Accuracy	F1
CAER-Net-S [15]	73.51	-
$FO_{R34}$	77.35 (0.002)	77.13 (0.002)
<b>MCF (<math>SAG-MHA_{enc}</math>)</b>	<b>79.63 (0.003)</b>	<b>79.36 (0.003)</b>

**Table 2.** Comparison of *MCF* module with state of the art in CAER-S (Test). **F1:** macro-F1, **FO:** Face only, **R34:** Resnet34 fully finetuned. Average of runs with 5 random seeds reported with standard deviation for *MCF* and Resnet34.

### 6.2. Ablation studies

We analyze the importance of different input context streams and associated models in EMOTIC. From Table 3, we can see that the Resnet34 model ( $PO_{R34}$ ) fully finetuned using person specific crops performs worst, thus indicating the need of additional contextual information. Inclusion of scene representation (cls token) from ViT pretrained on Places2 dataset with  $PO_{R34}$  via late fusion (LF) improves the mAP to **25.53**. Freezing  $PO_{R34}$  model followed by cross modal interaction through  $CM_{enc}$  composed of  $MHA_{enc}$  layers (visual-scene guided context stream in Fig. 3) further increases the mAP to **27.37**. Further we can see that fusion of both foreground and visual scene based context information through *MCF*



**Fig. 4.** Examples (a) and (b) from EMOTIC showing comparisons between top-5 predictions of  $PO_{R34}$  and MCF (MHA<sub>enc</sub>). Examples (c) and (d) from CAER-S showing comparisons between top predictions of  $FO_{R34}$  and MCF (SAG-MHA<sub>enc</sub>).  $PO$ : Person-only.  $FO$ : Face-only.  $R_{34}$ : Resnet34 fully finetuned. GT: Ground truth. Person or face instances marked by green bounding boxes

(MHA<sub>enc</sub>) results in the best performance (**29.53**). For both CM<sub>Enc</sub> (MHA<sub>enc</sub> + VS) and MCF (MHA<sub>enc</sub>), we freeze the Resnet34 model for extracting representations from person crops. For CAER-S, we

with hand enables MCF (SAG-MHA<sub>enc</sub>) to make a correct prediction (**Angry**) as compared to  $FO_{R34}$  (Resnet34 finetuned on face crops).

Model	mAP	$\lambda_1$	$\lambda_2$
PO <sub>R34</sub>	23.46 (0.006)	0.95	0.05
PO <sub>R34</sub> + VS + LF	25.53 (0.025)	0.6	0.4
CM <sub>Enc</sub> (MHA <sub>enc</sub> + VS)	27.37 (0.004)	0.5	0.5
<b>MCF (MHA<sub>enc</sub>)</b>	<b>29.53 (0.001)</b>	<b>0.8</b>	<b>0.2</b>

**Table 3.** Ablation study on different context streams and associated models for EMOTIC. **PO**: Person-only model, **R34**: Resnet34 fully fine-tuned, **VS**: Visual scene, **cls**: cls token, **LF**: Late fusion,  $\lambda_1$ : BCE weight,  $\lambda_2$ : MSE weight. Average of runs with 5 random seeds reported with standard deviation for the models

can see from Table 4 that inclusion of SAG-MHA<sub>enc</sub> layer instead of MHA<sub>enc</sub> improves the accuracy from **76.89** to **79.63**. This can be attributed to the self-attention based augmentation operation for the query features i.e. face representations from Resnet34 in SAG-MHA<sub>enc</sub> layer (Fig 2). For both MCF (MHA<sub>enc</sub>) and MCF (SAG-MHA<sub>enc</sub>), we finetune the Resnet34 model completely with MCF for extracting representations from face crops.

Model	Accuracy	F1
FO <sub>R34</sub>	77.35 (0.002)	77.13 (0.002)
<b>MCF (SAG-MHA<sub>enc</sub>)</b>	<b>79.63 (0.003)</b>	<b>79.36 (0.003)</b>
MCF (MHA <sub>enc</sub> )	76.89 (0.003)	76.74 (0.002)

**Table 4.** Ablation study on different context streams and associated models for CAER-S. **FO**: Face-only model, **R34**: Resnet-34 fully fine-tuned, Average of runs with 5 random seeds reported with standard deviation for the models

### 6.3. Qualitative examples

In Fig 4 (a) and (b), we can see that the inclusion of foreground context through captions like *basketball team celebrating* and *man crouching down with his hands on his face* results in consistent performance of MCF (MHA<sub>enc</sub>) as compared to  $PO_{R34}$  (Resnet34 finetuned on person crops). Similarly, for TV shows in Fig 4 (c) while the face crop based prediction from  $FO_{R34}$  is **sad**, inclusion of foreground context with visual scene information gives a correct prediction for MCF (SAG-MHA<sub>enc</sub>). In Fig 4 (d), the act of *gesturing*

## 7. CONCLUSION

In this work, we explore the role of contextual information in estimating human emotions with respect to the domains of natural scenes (EMOTIC) and TV shows (CAER-S). Since multimodal-VLN models are pretrained on large-scale image-text pairs from the web, we utilize their capabilities to obtain foreground context information in terms of descriptive captions. Further, we propose a purely attention-based multimodal context fusion (MCF) module to combine person-specific information with the visual scene and foreground context representations. Future work involves the extension of the MCF module to include geometric aspects of person context including pose information and evaluation using media-centered data like movies and advertisements.

## 8. ACKNOWLEDGEMENT

We would like to thank the Center for Computational Media Intelligence at USC for supporting this study.

## 9. REFERENCES

- [1] Daniel Dukes, Kathryn Abrams, Ralph Adolphs, Mohammed E Ahmed, Andrew Beatty, Kent C Berridge, et al., “The rise of affectivism,” *Nature human behaviour*, vol. 5, no. 7, pp. 816–820, 2021.
- [2] Sicheng Zhao, Xingxu Yao, Jufeng Yang, Guoli Jia, et al., “Affective image content analysis: Two decades review and new perspectives,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6729–6751, 2022.
- [3] S. Latif, R. Rana, S. Khalifa, R. Jurdak, J. Qadir, and B. W. Schuller, “Survey of deep representation learning for speech emotion recognition,” *IEEE Transactions on Affective Computing*, no. 01, pp. 1–1, sep 5555.
- [4] Chiara Zucco, Barbara Calabrese, and Mario Cannataro, “Sentiment analysis and affective computing for depression monitoring,” in *2017 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2017, pp. 1988–1995.



- [5] Tanaya Guha, Zhaojun Yang, Ruth B. Grossman, and Shrikanth S. Narayanan, "A computational study of expressive facial dynamics in children with autism," *IEEE Transactions on Affective Computing*, vol. 9, no. 1, pp. 14–20, 2018.
- [6] Andrey V. Savchenko, Lyudmila V. Savchenko, and Ilya Makarov, "Classifying emotions and engagement in online learning based on a single facial expression recognition neural network," *IEEE Transactions on Affective Computing*, pp. 1–12, 2022.
- [7] Xingxun Jiang, Yuan Zong, et al., "Dfew: A large-scale database for recognizing dynamic facial expressions in the wild," in *Proceedings of the 28th ACM International Conference on Multimedia*, New York, NY, USA, 2020, MM '20, p. 2881–2889, Association for Computing Machinery.
- [8] Ali Mollahosseini, Behzad Hasani, and Mohammad H. Mahoor, "Affectnet: A database for facial expression, valence, and arousal computing in the wild," *IEEE Transactions on Affective Computing*, vol. 10, pp. 18–31, 2019.
- [9] Carlos Busso, Zhigang Deng, Serdar Yildirim, Murtaza Bulut, and Chul Min others Lee, "Analysis of emotion recognition using facial expressions, speech and multimodal information," in *Proceedings of the 6th international conference on Multimodal interfaces*, 2004, pp. 205–211.
- [10] Lisa Feldman Barrett, Batja Mesquita, and Maria Gendron, "Context in emotion perception," *Current Directions in Psychological Science*, vol. 20, no. 5, pp. 286–290, 2011.
- [11] Moshe Bar, "Visual objects in context," *Nature Reviews Neuroscience*, vol. 5, pp. 617–629, 2004.
- [12] Peng Wang, An Yang, Rui Men, Junyang Lin, et al., "Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework," *CoRR*, vol. abs/2202.03052, 2022.
- [13] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal, "Unifying vision-and-language tasks via text generation," in *ICML*, 2021.
- [14] Ronak Kosti, Jose M. Alvarez, Adria Recasens, and Agata Lapedriza, "Context based emotion recognition using emotic dataset," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 11, pp. 2755–2766, 2020.
- [15] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and Kwanghoon Sohn, "Context-aware emotion recognition networks," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 10142–10151.
- [16] Takahiko Masuda, Phoebe C. Ellsworth, Batja Mesquita, et al., "Placing the face in context: cultural differences in the perception of facial emotion.," *Journal of personality and social psychology*, vol. 94 3, pp. 365–81, 2008.
- [17] Bernd Dudzik, Michel-Pierre Jansen, Franziska Burger, Frank Kaptein, et al., "Context in human emotion perception for automatic affect detection: A survey of audiovisual databases," *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pp. 206–212, 2019.
- [18] Matthias J. Wieser and Tobias Brosch, "Faces in context: A review and systematization of contextual influences on affective face processing," *Frontiers in Psychology*, vol. 3, 2012.
- [19] Ali Mollahosseini, Behzad Hasani, and Mohammad H. Mahoor, "Affectnet: A database for facial expression, valence, and arousal computing in the wild," *IEEE Transactions on Affective Computing*, vol. 10, no. 1, pp. 18–31, 2019.
- [20] Emad Barsoum, Cha Zhang, Cristian Canton Ferrer, and Zhengyou Zhang, "Training deep networks for facial expression recognition with crowd-sourced label distribution," in *ACM International Conference on Multimodal Interaction (ICMI)*, 2016.
- [21] Minghui Zhang, Yumeng Liang, and Huadong Ma, "Context-aware affective graph reasoning for emotion recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*, 2019, pp. 151–156.
- [22] Trisha Mittal, Pooja Guhan, et al., "Emoticon: Context-aware multimodal emotion recognition using frege's principle," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14222–14231, 2020.
- [23] Ioannis Pikoulis, Panagiotis P. Filntisis, and Petros Maragos, "Leveraging semantic scene characteristics and multi-stream convolutional architectures in a contextual approach for video-based visual emotion recognition in the wild," in *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*, 2021.
- [24] Junnan Li, Ramprasaath R. Selvaraju, et al., "Align before fuse: Vision and language representation learning with momentum distillation," in *NeurIPS*, 2021.
- [25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *ArXiv*, vol. abs/2010.11929, 2021.
- [26] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba, "Places: A 10 million image database for scene recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [27] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, 2019.
- [29] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, et al., "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Online, Oct. 2020, pp. 38–45, Association for Computational Linguistics.
- [30] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian, "Deep modular co-attention networks for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6281–6290.
- [31] Ilya Loshchilov and Frank Hutter, "Decoupled weight decay regularization," in *ICLR*, 2019.
- [32] Diederik P. Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2015.
- [33] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, pp. 1499–1503, 2016.
- [34] Adam Paszke, Sam Gross, Francisco Massa, et al., "Pytorch: An imperative style, high-performance deep learning library," *ArXiv*, vol. abs/1912.01703, 2019.