

# Stimuli-Aware Visual Emotion Analysis

Jingyuan Yang, Jie Li, Xiumei Wang, Yuxuan Ding and Xinbo Gao, *Senior Member, IEEE*

**Abstract**—Visual emotion analysis (VEA) has attracted great attention recently, due to the increasing tendency of expressing and understanding emotions through images on social networks. Different from traditional vision tasks, VEA is inherently more challenging since it involves a much higher level of complexity and ambiguity in human cognitive process. Most of the existing methods adopt deep learning techniques to extract general features from the whole image, disregarding the specific features evoked by various emotional stimuli. Inspired by the *Stimuli-Organism-Response (S-O-R)* emotion model in psychological theory, we proposed a stimuli-aware VEA method consisting of three stages, namely stimuli selection (S), feature extraction (O) and emotion prediction (R). First, specific emotional stimuli (*i.e.*, color, object, face) are selected from images by employing the off-the-shelf tools. To the best of our knowledge, it is the first time to introduce stimuli selection process into VEA in an end-to-end network. Then, we design three specific networks, *i.e.*, Global-Net, Semantic-Net and Expression-Net, to extract distinct emotional features from different stimuli simultaneously. Finally, benefiting from the inherent structure of Mikel's wheel, we design a novel hierarchical cross-entropy loss to distinguish hard false examples from easy ones in an emotion-specific manner. Experiments demonstrate that the proposed method consistently outperforms the state-of-the-art approaches on four public visual emotion datasets. Ablation study and visualizations further prove the validity and interpretability of our method.

**Index Terms**—Visual emotion analysis, emotional stimuli, hierarchical cross-entropy loss, emotion classification, convolutional neural networks.

## I. INTRODUCTION

In recent years, after the breathtaking success in traditional computer vision tasks such as object detection [1]–[3] and image classification [4]–[6], researchers have gradually turned their eyes from low-level vision tasks to high-level ones, including visual reasoning [7], image aesthetic assessment [8], visual emotion analysis [9], *etc.* Among all the high-level vision tasks, *Visual Emotion Analysis (VEA)* is one of the most challenging tasks for the existing *affective gap* [10] between low-level pixels and high-level emotions. Against all odds, VEA is still promising as understanding human emotions is a crucial step towards strong artificial intelligence [11]. As



Fig. 1. Mikel's wheel from psychological model. Affective images from FI dataset with emotion categories (*i.e.*, amusement, anger, awe, contentment, disgust, excitement, fear, and sad) and polarities (*i.e.*, positive, negative).

shown in Fig. 1, people may experience different emotions towards different images according to Mikel's wheel [12], [13]. Consequently, VEA aims at finding out how people feel emotionally towards different visual information, which has become an important research topic recently. Solving the VEA task may have a tremendous impact on real-world applications such as opinion mining [14], [15], smart advertising [16], [17], and mental disease treatment [18], [19].

Various methods have been proposed to deal with this challenging yet promising problem, including earlier traditional methods and recent deep learning ones. In an early stage, inspired by psychological and aesthetic theories, researchers designed assorted hand-crafted features manually, which included color, texture, composition, balance, emphasis, *etc.* [20]–[24]. Most of the early attempts designed specific yet limited features, which failed to cover all important emotional factors and thus resulted in degraded performance on large-scale datasets. With the rapid development of *Convolutional Neural Networks (CNNs)*, more and more researchers employed deep learning networks to VEA [9], [25]–[30]. Instead of designing emotion features manually, CNNs were capable of mining emotional representations automatically in an end-to-end manner and consequently achieved better results. Nevertheless, most of the existing deep learning methods directly extracted general features from the whole image, which failed to consider the unique evocation process involved in VEA.

In view of the significance of VEA, great attention has been paid to it not only in computer vision but also in psychology. Psychologist Frijda [31] found that some special substances evoked emotion and named them emotional

Manuscript received January 21, 2020; revised April 6, 2021 and May 11, 2021; accepted August 9, 2021. This work was supported in part by the National Natural Science Foundation of China under Grants 62036007, 62050175 and 61772402. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Trac D. Tran. (*Corresponding author: Xinbo Gao.*)

J. Yang, J. Li, X. Wang and Y. Ding are with the School of Electronic Engineering, Xidian University, Xi'an 710071, China (e-mail: jingyuanyang@stu.xidian.edu.cn; leejie@mail.xidian.edu.cn; wangxm@xidian.edu.cn; yxding@stu.xidian.edu.cn).

X. Gao is with the School of Electronic Engineering, Xidian University, Xi'an 710071, China (e-mail: xbgao@mail.xidian.edu.cn) and with the Chongqing Key Laboratory of Image Cognition, Chongqing University of Posts and Telecommunications, Chongqing 400065, China (e-mail: gaobx@cqupt.edu.cn).



Fig. 2. Examples from FI dataset. Different emotional stimuli (*i.e.*, color, object, face) may evoke different emotions (*i.e.*, positive, negative).

stimuli. Psychologists Mehrabian and Russell [32] suggested that people perceived emotions through three steps – *stimuli*, *organism* and *response*, which was the so-called *S-O-R* model. In traditional computer vision tasks, networks often consist of two parts, namely feature extraction and prediction, which can be regarded as organism (O) and response (R) in S-O-R model, correspondingly. However, S-O-R model suggests that it is specific stimuli rather than the whole image that evoke emotions, which makes stimuli selection (S) indispensable to VEA compared with other vision tasks. Therefore, as opposed to extracting features from the whole image directly, we first select emotional stimuli from the image and extract specific features from those stimuli afterwards.

As mentioned in a psychological literature [33], color, certain objects, facial expressions or any other attributes can be summarized as emotional stimuli. When it comes to a scenery image without any salient objects, the emotion is evoked by global features, such as color, brightness, texture, *etc.* When several objects are presented in an image, people tend to deduce the semantic correlations between those objects and subsequently generate an emotion. When human face appears in an image, we may focus on the facial expression and evoke a similar emotion as well, which is called empathy [34]. As shown in Fig. 2, under different emotional stimuli (*i.e.*, color, object, and face), people may experience different emotions. For example, sunset may bring people a negative feeling while smiling faces often deliver a positive emotion. Therefore, in the proposed method, we choose color, object and face as special emotional stimuli both theoretically [33] and empirically.

Therefore, motivated by the above facts from psychology, we proposed a novel stimuli-aware VEA network consisting of three stages, namely stimuli selection (S), feature extraction (O) and emotion prediction (R). In the first stage, special emotional stimuli (*i.e.*, color, object and face) are selected by employing the off-the-shelf tools in deep learning. After that, we design three specific networks to extract distinct features from different stimuli simultaneously, namely Global-Net, Semantic-Net and Expression-Net. Eventually, we leverage all the above features for emotion prediction.

Categories in most classification tasks are irrelevant to each other, while there exists an inherent hierarchical structure between emotions according to Mikel’s wheel [12], [13], a widely-used emotion model from psychology. As shown in Fig. 1, emotions on Mikel’s wheel are separated into eight categories (*i.e.*, *amusement*, *anger*, *awe*, *contentment*,

*disgust*, *excitement*, *fear*, and *sad*) as well as two polarities (*i.e.*, *positive* and *negative*). To be specific, amusement, awe, contentment, excitement belong to positive emotions while anger, disgust, fear, sad are negative ones, where different colors (*i.e.*, red and blue) in Fig. 1 denote positive and negative emotions respectively. To exploit such prior knowledge, we propose a novel hierarchical cross-entropy loss for VEA. In general classification tasks, examples can be divided into true examples and false ones considering whether their emotion categories are correctly classified. We further distinguish hard false examples from easy false ones by judging whether their emotion polarities are correctly classified. To be specific, we add an auxiliary polarity loss to the traditional emotion loss, with increased penalties on hard false examples compared with easy ones. By implementing the proposed hierarchical cross-entropy loss, the whole network can pay more attention to hard false examples and consequently predict emotions in a hierarchical and emotion-specific manner.

Our contributions can be summarized as follows:

- We propose a stimuli-aware VEA network consisting of three stages, namely stimuli selection (S), feature extraction (O) and emotion prediction (R), which corresponds to the three stages in S-O-R model from psychology. The proposed method consistently outperforms the state-of-the-art methods on four public visual emotion datasets.
- We choose color, object and face as special emotional stimuli theoretically and empirically, and design three specific networks to extract features from different stimuli simultaneously. To the best of our knowledge, it is the first work that predict emotions from specific emotional stimuli instead of the whole image.
- We design a novel hierarchical cross-entropy loss to distinguish hard false examples from easy ones with increased penalties, which benefits from the inherent structure of psychological model and further boosts the classification result in an emotion-specific manner.

The rest of the paper is organized as follows. Section II overviews the existing methods on visual emotion analysis, object detection and facial expression recognition. In Section III, we introduce our proposed stimuli-aware VEA network and the novel hierarchical cross-entropy loss. Extensive experiments and visualization results on four public datasets are given in Section IV and Section V. Finally, we conclude our work in Section VI.

## II. RELATED WORK

In this section, we review the previous methods on VEA from earlier traditional methods and recent deep learning ones. Meanwhile, related work on object detection and facial expression recognition is also included in this section.

### A. Visual Emotion Analysis

Emotion psychologists established two models describing emotions, namely Categorical Emotion States (CES), and Dimensional Emotion Space (DES). CES model divided emotions into several categories, which mainly included a two-category approach and an eight-category approach. DES model

employed a continuous Valence-Arousal-Dominance (VAD) space to separate different emotions. As for simplicity and intuitiveness, most researchers preferred CES model to DES model. Some typical models in CES included Ekman's six emotion categories [35] and Mikel's eight emotion categories [12], [13], serving as our basic models as well.

1) *Traditional Methods*: Earlier work on VEA mainly focused on the design of traditional hand-crafted features. Machajdik *et al.* [20] extracted low-level features (*i.e.*, color, texture, composition, and content), and combined them to predict emotions of an image. Siersdorfer *et al.* [22] adopted global and local RGB histogram, as well as SIFT-based bag of visual terms as emotion features. Zhao *et al.* [21] investigated the relationship between artistic principles and emotions by extracting principle-of-art-based emotion features, which consisted of balance, emphasis, harmony, variety, gradation, and movement. Borth *et al.* [24] proposed a concept named Adjective Noun Pairs (ANPs) to simultaneously preserve emotional and location information of objects in an image. In order to sum up different level of emotion features, low-level features from elements-of-art, mid-level features from principles-of-art, as well as high-level features from ANPs and facial expressions were combined in a multi-graph learning manner by Zhao *et al.* [23]. Chen *et al.* [24] implemented object detection methods to recognize the top six frequent objects (*i.e.*, car, dog, dress, face, flower, and food), and proposed a novel classification approach to model the concept similarity between ANPs. These hand-crafted features were specifically designed yet limited to represent, which failed to cover all important factors in VEA and resulted in degraded performance on large-scale datasets.

2) *Deep Learning Methods*: With the rapid development of deep learning, *Convolutional Neural Networks (CNNs)* were applied in various computer vision tasks and achieved the state-of-the-art performance. Due to its great success, researchers in VEA also exploited deep CNN for sentiment prediction. Based on their previous work SentiBank [24], Chen *et al.* constructed DeepSentiBank [25] by replacing the SVM classifier with deep CNN, which resulted in a great improvement in both annotation accuracy and retrieval performance. You *et al.* [30] proposed a novel progressive CNN architecture (PCNN) to make use of the noisy labeled data for binary sentiment classification. To help training with deep CNN, You *et al.* [36] established a benchmark for emotion recognition by designing a large scale dataset based on the Mikel's eight emotion categories. A multi-level deep representation network (MldrNet) [26] was proposed by Rao *et al.* to extract emotion features from multi-scale patches (*i.e.*, pixel-level feature, aesthetic feature and semantic feature), which was an early attempt to design an emotion-specific network. In order to integrate the features extracted from different levels more effectively, Zhu *et al.* [37] adopted Bi-GRU as feature fusion module while keep traditional CNN as feature extraction module. Further, Rao *et al.* [38] proposed a multi-level deep network to utilize both low-level and high-level features for predicting visual emotions. Recently, Yang *et al.* [27] constructed a global-local network, implemented object detection method to find affective regions in local

branch and combined two branches to make final sentiment prediction. Yang *et al.* [9] proposed a weakly supervised coupled network (WSCNet) with two branches to leverage the localized information through attention mechanism. Similarly, He *et al.* [39] proposed a multi-attentive pyramidal mode, aiming to find emotional cues from local regions and their relationships. A novel CNN model was proposed by Zhang *et al.* [40] to extract and integrate content information as well as style information to infer visual emotions. From a higher perspective, Zhang *et al.* [41] propose a novel object semantics sentiment correlation model (OSSCM) to predict emotions from object semantics. However, most of the existing deep learning methods directly fed the whole image into general deep networks for feature extraction, without a careful consideration of the unique evocation process involved in VEA.

Traditional methods designed specific emotion features manually, which resulted in limited representation ability and degraded performance when facing massive data. Oppositely, deep learning methods were capable of extracting effective features automatically yet lacked interpretability. Learning from both approaches, we first select emotional stimuli manually based on the psychological theories and then implement deep networks to extract emotion features from selected stimuli. Based on S-O-R model, we construct a manually stimuli selection and a deep feature extraction network, with both advantages of interpretability and effectiveness. Besides, compared with previous methods, we make the first attempt to introduce stimuli selection into VEA, which is proved interpretable and effective.

## B. Object Detection

As one of the basic problems in computer vision, object detection was heated discussed and consequently made great progress recently. Due to its high accuracy and efficiency, object detection was implemented into various tasks as a preprocessing stage, such as image caption, VQA, and visual reasoning [42]. Among all the object detection algorithms, several benchmarks were worth-mentioning, including R-CNN [1], fast R-CNN [2], and faster R-CNN [3]. R-CNN, as a pioneering work in this field, implemented deep CNN into object detection for the first time, which greatly improved the detection performance. By replacing the original serial structure with a parallel structure, Girshick proposed a new method and named it fast R-CNN, which led to a higher speed as well as improved accuracy. The biggest innovation in faster R-CNN was the Region Proposal Network (RPN), as opposed to the selective-search in previous methods. RPN not only made it possible to generate effective proposals in an end-to-end network, but it also improved the detection speed to a large extent. Considering its great advantage in both speed and accuracy, we employ faster R-CNN as the object detection method in our network.

## C. Facial Expression Recognition

*Facial expression recognition (FER)* aims to recognize facial changes in response to a person's internal emotional states, which drawn great interest among researchers and was widely

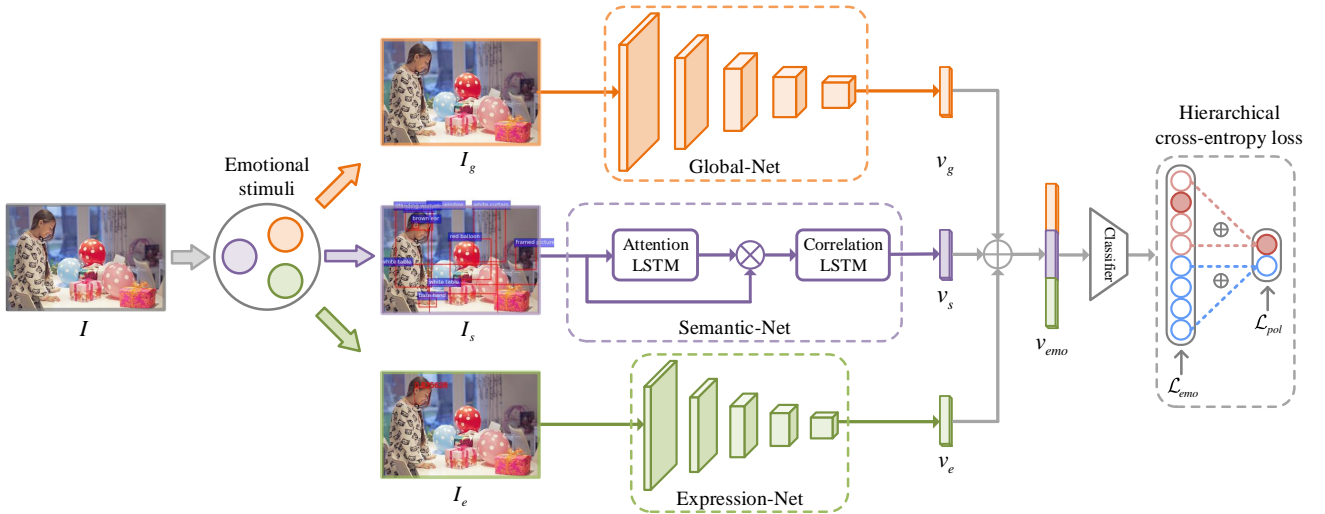


Fig. 3. Framework of the proposed stimuli-aware VEA network. Based on the S-O-R psychological model, our network consists of three stages, namely stimuli selection, feature extraction and sentiment prediction. First, specific emotional stimuli (*i.e.*, color, object, face) are selected from images by employing the off-the-shelf tools in deep learning methods (Sec. III-A). Then, we construct three specific networks (*i.e.*, Global-Net, Semantic-Net and Expression-Net) to extract distinct emotion features from different stimuli simultaneously (Sec. III-B). Finally, considering the inherent structure of emotion categories in Mikel’s wheel, we propose a novel hierarchical cross-entropy loss penalizing hard false examples from easy ones (Sec. III-C).

used in real-world applications in recent years [43]–[48]. In FER, it is essential to pre-process the image with face alignment, data augmentation and face normalization [49]. In face alignment, the first step is to detect the face [49]. Dlib [50] and MTCNN [51] are two of the most commonly-used face detectors in FER. Considering its wide applications in recent years [52], [53], we employ Dlib as our face detector. After pre-processing, the face image is then sent into a feature extraction network, in which CNN is a common practice. Various datasets were collected for facial expression recognition, including lab datasets and web datasets. Lab datasets consisted of face images specially take in the laboratory environment [54]–[56], while web datasets were gathered from the internet [57]–[59]. As most of our emotion datasets were collected from the web, we initialize our Expression-Net with pre-trained parameters from a widely used web dataset named FER2013 [57]. FER2013 was a large-scale dataset collected automatically on the web by Google image search API, which was introduced during the ICML 2013 Challenges in Representation Learning.

### III. METHODOLOGY

In this section, we propose a novel stimuli-aware network for emotion predictions through mining emotions from different emotional stimuli. Inspired by the S-O-R psychological model, our network is designed with three stages, *i.e.*, stimuli selection (S), feature extraction (O) and emotion prediction (R). To the best of our knowledge, it is the first time to introduce stimuli selection (S) into VEA in an end-to-end network. Fig. 3 shows the architecture of the proposed method. Firstly, specific emotional stimuli (*i.e.*, color, object, face) are selected from images by employing the off-the-shelf detection tools in deep learning methods (Sec. III-A). Considering the uniqueness of different stimuli, we construct three specific

branches, namely Global-Net, Semantic-Net and Expression-Net, to extract distinct emotional features from different stimuli simultaneously (Sec. III-B). Finally, benefiting from the inherent structure of emotion categories, we design a hierarchical cross-entropy loss to distinguish hard false examples from easy ones with an increased penalty (Sec. III-C).

#### A. Stimuli Selection

Unlike other traditional computer vision tasks [1]–[6], VEA is much more challenging as emotions are complex and ambiguous, which cannot be “found” directly in an image. Apparently, it is easy to classify an image to the category of dog while it is thorny to judge an image with sad. Therefore, how to effectively bridge the *affective gap* [10] between low-level pixels and high-level emotions has become the biggest challenge. Fortunately, great attention has been paid to VEA not only in computer vision but also in psychology. Psychologist Frijda [31] found out that some special substances may evoke emotion and named them emotional stimuli. Different from the previous stimuli-response (S-R) model, psychologist Woodworth [60] argued the importance of organism (O). Based on this, psychologists Mehrabian and Russell [32] further suggested that human perceive emotions through three steps – stimuli, organism and response, which was the so-called S-O-R model [61]–[63]. As mentioned in [33], “How we perceive our environment is thus shaped by categorization processes which guide and constrain the organization of incoming stimulus information,” indicating the importance of stimuli (S) in S-O-R model. Further, psychologist Brosch suggested that color, certain objects, facial expressions or any other attributes can be categorized as emotional stimuli [33]. Therefore, unlike previous methods, we introduce stimuli selection to VEA, aiming to predict emotions from specific emotional stimuli instead of the whole image. Inspired by the S-O-R model, we construct our network with three stages, *i.e.*, stimuli selection



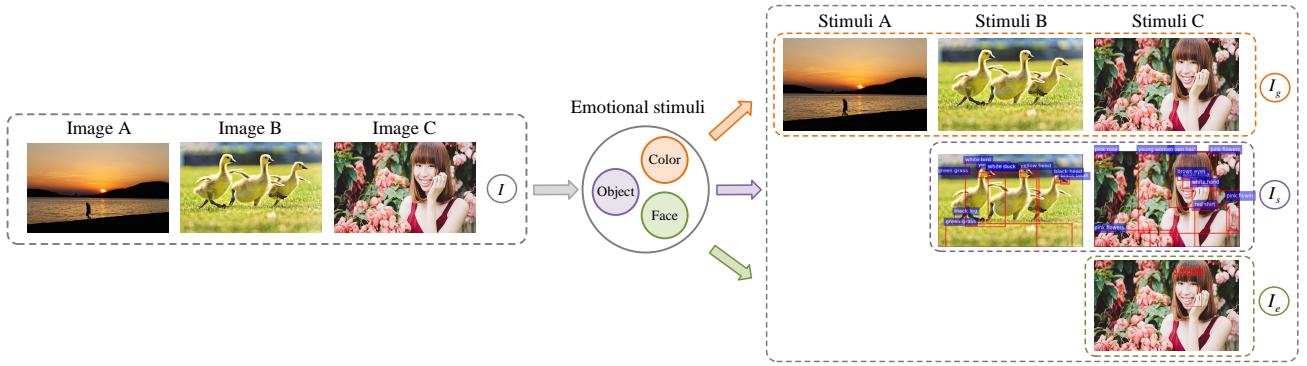


Fig. 4. Overview of the stimuli selection process. Based on psychological theory, we choose color, object and face as typical emotional stimuli and implement the off-the-shelf deep learning tools to select these stimuli from our datasets. Through this process, an affective image is turned into a set of emotional stimuli, which are specifically designed and further utilized for emotion prediction.

(S), feature extraction (O) and emotion prediction (R), where color, object and face are chosen as typical emotional stimuli according to psychological theories.

We further illustrate the stimuli selection process with three examples in Fig. 4. Image A is a scenery image describing sunset, with hardly any object or face in it. After fed into the stimuli selection module, stimuli A is generated from image A with color stimulus alone. We can infer from the pale-yellow color stimulus that image A implies *sad* emotion. Image B shows three little cute ducks running on the grass. Different from the previous scenery image, there are distinct objects in image B, which can be viewed as object stimulus. Both object stimulus and color stimulus contribute to the final emotion prediction of *awe*. Image C describes a smiling girl standing behind a flower-covered wall, which seems comfortable and satisfied. Feeding image C into the stimuli selection module and we will get stimuli C containing face stimulus, object stimulus and color stimulus simultaneously. Specifically, we can see a happy girl in face stimulus, a wall covered with pink flowers in object stimulus, and the color of the entire picture in color stimulus. All the three emotional stimuli are indispensable to the final emotion prediction of *contentment*. To sum up, different images may have different combinations of emotional stimuli. Once any of them occurred, it contributes to the final emotion prediction, alone or together with other emotional stimulus.

Therefore, we construct three branches for stimuli selection, *i.e.*, color, object and face. To be specific, we keep the entire image for color stimulus, *i.e.*,  $I_g = I$ . Meanwhile, object stimulus ( $I_s$ ) and face stimulus ( $I_e$ ) are further selected from the entire image by employing the off-the-shelf detection tools in deep learning methods.

For the object branch, we employ faster R-CNN as our stimulus detector. To obtain a concrete and precise description of objects, we initialize our faster R-CNN with the pre-trained parameters from Visual Genome dataset [64], replacing the original PASCAL VOC [65] and MS COCO [66] datasets. Specifically, in the training process, our object detector is only feed forward with a set of frozen parameters. In other words, our object detector is only a general detection tool which is never specially optimized using emotion datasets alone. We

visualize the object proposals in stimuli selection as is shown in Fig. 4. As discussed above, strong correspondence does not exist between a single object and a certain emotion, which makes it hard to predict the emotion of the whole image from a single object. Hence, we select multiple objects from an image by ranking their detection confidence. Specifically, after stimuli selection, image  $I$  is turned into object stimulus  $I_s$ , containing top- $N$  detected objects. We denote the object proposals as  $I_s = \{i_1, \dots, i_N\}$  and the object features for each proposal as  $\mathbf{F}_s = \{\mathbf{f}_1, \dots, \mathbf{f}_N\}$ , which are automatically extracted from faster-RCNN.  $N$  denotes the number of regions with  $N = 10$ , which is further ablated with experiments in Sec. IV-D2. Considering the possible redundancy caused by multiple objects, a specific feature extraction network is further proposed to infer the semantic correlations between different object features  $\mathbf{F}_s$ , which will be discussed in Sec. III-B2.

In order to select face stimulus, we employ landmark face detector Dlib [50] from facial expression recognition. In detail, Dlib first detects faces in each image and then crops them into an aligned size automatically. To be specific, taking the  $448 \times 448$  image  $I$  as input, the face detector will output a face crop  $I_e$  with size  $48 \times 48$ , which can be regarded as the face stimulus. It is obvious that not all the images in FI dataset contain human faces. If an image does not contain any face, the face detector will output nothing. Besides, when multiple faces appear in an image at the same time, we only choose the largest face as  $I_e$ , assuming that people are mostly likely to be attracted and thus affected by it. Through the stimuli selection process, an affective image is turned into a set of emotional stimuli, containing color, object and face, which are specifically designed for further emotion prediction in the following sections.

### B. Feature Extraction

After stimuli selection,  $I_g$ ,  $I_s$  and  $I_e$  are generated as emotional stimuli representing image  $I$ . Furthermore, we construct three specific branches, namely Global-Net, Semantic-Net and Expression-Net, to extract distinct emotion features from different stimuli simultaneously.

1) *Global-Net*: As mentioned in Sec. III-A, color, objects, faces or other attributes can be summarized as emotional

stimuli. Color feature is a global feature showing a strong correlation with emotions. However, other global features (*e.g.*, texture, pattern, composition, scene) also contribute to the evocation process of emotion and thus cannot be neglected. Besides, in deep learning methods, there is no specialized network for color feature extraction. Therefore, we construct Global-Net to extract color feature and other global features.

Most of the existing methods [9], [27] in VEA employ Convolutional Neural Networks (CNNs) to extract global features, due to the powerful representation ability of deep features. In this work, our Global-Net is based on ResNet-50 [6], balancing both efficiency and accuracy. Specifically, Global-Net consists of five convolutional layers and a Global Averaging Pooling (GAP) layer, which is described in Eq. (1). In Global-Net, global stimulus  $I_g$  is first sent into the fully convolutional networks of ResNet-50 and then fed into the GAP layer to obtain the global feature vector:

$$\mathbf{F}_g = \text{FCN}_{\text{res50}}(I_g), \quad (1)$$

$$\mathbf{v}_g = G_{\text{avg}}(\mathbf{F}_g), \quad (2)$$

where  $\text{FCN}_{\text{res50}}$  is the fully convolutional networks in ResNet-50 and  $G_{\text{avg}}$  is the following GAP layer.  $\mathbf{F}_g \in \mathbb{R}^{d_1 \times w \times h}$  denotes the feature maps extracting from the last convolutional layer and  $\mathbf{v}_g \in \mathbb{R}^{d_1}$  represents the extracted global vector. Specifically,  $w, h$  are the spatial size of the feature map while  $d_1$  is the length of the global vector with  $d_1 = 2048$ .

2) *Semantic-Net*: In Sec. III-A, we implement object detection methods to select multiple objects in an image by ranking their detection confidence. Due to the complexity in emotion evocation process, it is not always an isolated object itself evokes a specific emotion, but the semantic correlations between different objects jointly determine the final result. For example, when we see white roses in Fig. 5(b), it is hard to decide whether a positive emotion or a negative one. However, as shown in Fig. 5(a), when white roses are accompanied by a bride, we can infer that it is a wedding and are very likely to evoke a positive emotion of *awe*. Oppositely, when a white rose and a tombstone come together in Fig. 5(c), we may have a negative feeling of *sad*. Therefore, instead of using the multi-object features directly, it is more reasonable to mine the semantic correlations between different objects.

Researchers have recently drawn great attention on exploiting the relationships between different objects, including Visual Question Answering (VQA), image captioning, visual reasoning, *etc.* In image caption, a newly released up-down attention model [67] has achieved great performance by implementing attention mechanism to weigh different object features. Based on this model, we design a emotion-specific Semantic-Net to mine the semantic correlations between different objects in our network. In order to exploit the correlations between different objects, we treat object features as a sequential data, and propose a Long-Short Term Memory (LSTM) [68] network to model such dependencies. The LSTM network is composed of two LSTM layers, namely attention LSTM and correlation LSTM, which is a standard implementation [69]. Firstly, we design an attention LSTM aiming to weigh each object features as well as reducing the redundancy of multiple

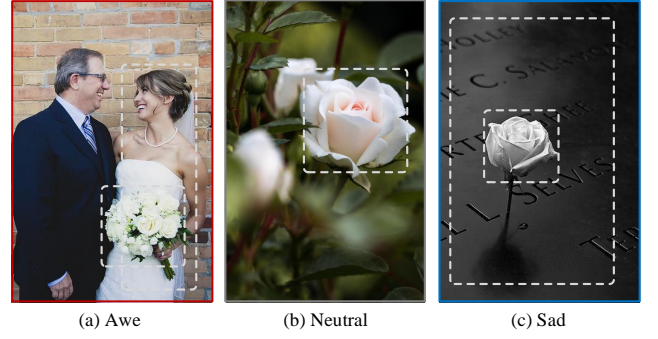


Fig. 5. Examples from FI dataset (*i.e.*, (a), (c)) and the internet (*i.e.*, (b)). It is not the white rose alone that evokes a specific emotion, but the interactions between different objects jointly determine the final result.

object features. After that, correlation LSTM is constructed to mine the semantic correlations between different objects.

After object stimulus selection, we obtain object proposals  $I_s = \{i_1, \dots, i_N\}$  and corresponding object features  $\mathbf{F}_s = \{\mathbf{f}_1, \dots, \mathbf{f}_N\}$ , in which  $N$  represents the number of objects in an image. In attention LSTM,  $\mathbf{x}_t^{\text{att}}$  denotes the input vector while  $\mathbf{h}_t^{\text{att}}$  denotes the output vector. The input vector  $\mathbf{x}_t^{\text{att}}$  to the attention LSTM at each time step consists of the previous hidden state  $\mathbf{h}_{t-1}^{\text{cor}}$  of the correlation LSTM, concatenated with the mean-pooled object features  $\mathbf{f}_{\text{mean}} = \frac{1}{N} \sum_{i=1}^N \mathbf{f}_i$ . The whole process of attention LSTM is given by

$$\mathbf{h}_t^{\text{att}} = \text{LSTM}_{\text{att}}(\mathbf{x}_t^{\text{att}}, \mathbf{h}_{t-1}^{\text{att}}), \quad (3)$$

$$\mathbf{x}_t^{\text{att}} = [\mathbf{h}_{t-1}^{\text{cor}}, \mathbf{f}_{\text{mean}}]. \quad (4)$$

The output  $\mathbf{h}_t^{\text{att}}$  of attention LSTM is then fed into Bahdanau attention [70] module, together with object features  $\mathbf{f}_i$ . The attention module generates a normalized attention weight  $\alpha_{i,t}$  for each object features  $\mathbf{f}_i$  at time step  $t$  as

$$a_{i,t} = \omega_a \tanh(\mathbf{W}_f \mathbf{f}_i + \mathbf{W}_h \mathbf{h}_t^{\text{att}}), \quad (5)$$

$$\alpha_{i,t} = \text{softmax}(a_{i,t}), \quad (6)$$

where  $\mathbf{W}_f \in \mathbb{R}^{M \times F}$ ,  $\mathbf{W}_h \in \mathbb{R}^{M \times H}$  and  $\omega_a \in \mathbb{R}^M$  are learned parameters of the attention module. The input object features  $\mathbf{F}_s = \{\mathbf{f}_1, \dots, \mathbf{f}_N\}$  are then weighed by the attention coefficient  $\alpha_{i,t}$ . The weighted object features, as shown in Eq. (7), are then fed into correlation LSTM:

$$\mathbf{f}_{\text{att}} = \sum_{i=1}^N \alpha_{i,t} \mathbf{f}_i. \quad (7)$$

In correlation LSTM,  $\mathbf{x}_t^{\text{cor}}$  denotes the input vector while  $\mathbf{h}_t^{\text{cor}}$  denotes the output vector. The input vector is concatenated by previous hidden state  $\mathbf{h}_{t-1}^{\text{att}}$  of the attention LSTM and the weighted object features  $\mathbf{f}_{\text{att}}$ , given by

$$\mathbf{h}_t^{\text{cor}} = \text{LSTM}_{\text{cor}}(\mathbf{x}_t^{\text{cor}}, \mathbf{h}_{t-1}^{\text{cor}}), \quad (8)$$

$$\mathbf{x}_t^{\text{cor}} = [\mathbf{h}_{t-1}^{\text{att}}, \mathbf{f}_{\text{att}}]. \quad (9)$$

After conducting both attention LSTM and correlation LSTM, not only has the redundancy between multiple objects been reduced, but the correlation between different objects has also been mined. Therefore, the output of the Semantic-Net can be

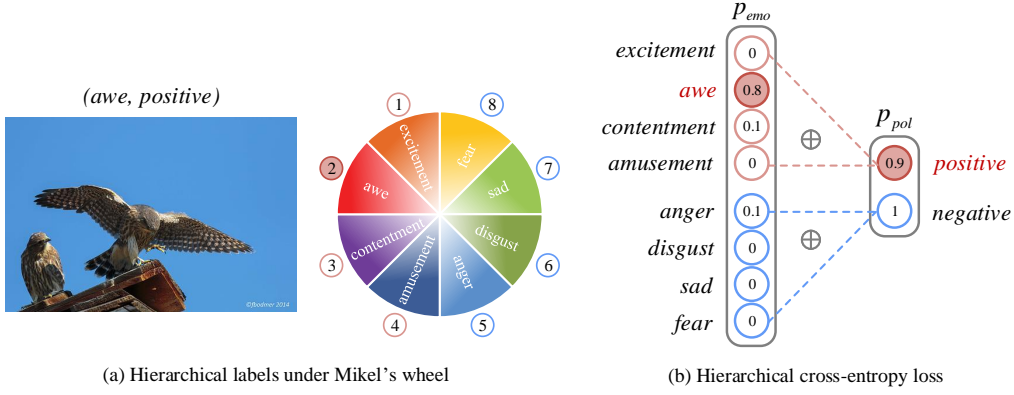


Fig. 6. Detailed illustration of the proposed hierarchical cross-entropy loss. Our hierarchical cross-entropy is motivated by the inherent structure of Mikel's wheel from psychology (a) and is further calculated corresponding to such hierarchical structure (b).

viewed as semantic feature with a higher-level representation of objects, which is given as

$$\mathbf{v}_s = \mathbf{h}_T^{cor}, \quad (10)$$

where  $T$  denotes the last time step in LSTM and  $\mathbf{v}_s \in \mathbb{R}^{d_2}$  represents the extracted semantic feature with  $d_2 = 512$ .

3) *Expression-Net*: It is obvious that human facial expression is the most direct way to convey emotions. In our daily life, if you see someone with a certain emotion, you will probably be affected by the emotion and evoke a similar emotion as well. This is called empathy, which was introduced by psychologist Jim Rogers [34]. Empathy is the ability to experience the inner world of other people, which is a peculiar ability of human beings. Therefore, facial expression is considered as an indispensable factor in VEA, which is designed as a branch in our network as well.

Since face images are small in size, previous researchers usually implement a shallow deep network to extract facial expression features. In light of this, we employ ResNet-18 as the base net to construct Expression-Net. Our Expression-Net is first initialized with pre-trained parameters from FER2013 dataset [57] and then jointly fine-tuned on FI dataset [36] with the rest networks. Expression-Net consists of five convolutional layers and a GAP layer, which is similar to Global-Net. It is worth mentioning that not all the images in our affective datasets have face stimulus. Therefore, the design of Expression-Net is divided into two cases as shown in Eq. (11). If face stimulus  $I_e$  exists, we first send it into the fully convolutional networks of ResNet-18 and then feed it into the GAP layer to obtain the expression feature  $\mathbf{v}_e$ . If face detector cannot find any face stimulus, we directly set  $\mathbf{v}_e$  to a zero vector, which does not contribute to the emotion prediction:

$$\mathbf{v}_e = \begin{cases} G_{avg}(\text{FCN}_{\text{res18}}(I_e)), \mathbf{v}_e \in \mathbb{R}^{d_3}, I_e \text{ exists} \\ \mathbf{0}, \text{ else} \end{cases} \quad (11)$$

where  $\text{FCN}_{\text{res18}}$  is the fully convolutional networks in ResNet-18 and  $G_{avg}$  is the following GAP layer. Besides,  $\mathbf{v}_e \in \mathbb{R}^{d_3}$  represents the extracted expression feature with  $d_3 = 512$ .

### C. Emotion Prediction

1) *Feature Fusion*: In the previous sections, we select three typical emotional stimuli (Sec. III-A) and design specific

networks to extract distinct features from them (Sec. III-B). As all the three features, *i.e.*, global feature ( $\mathbf{v}_g$ ), semantic feature ( $\mathbf{v}_s$ ), and expression feature ( $\mathbf{v}_e$ ), are indispensable and complementary in emotion evocation process, we further concatenate them as

$$\mathbf{v}_{emo} = \text{concat}[\mathbf{v}_g, \mathbf{v}_s, \mathbf{v}_e], \quad (12)$$

where  $\mathbf{v}_{emo} \in \mathbb{R}^{d_1+d_2+d_3}$  denotes emotion feature. The emotion feature  $\mathbf{v}_{emo}$  is further used for emotion prediction in Eq. (13).

2) *Hierarchical Cross-Entropy Loss*: In traditional classification tasks, cross-entropy (CE) loss has been widely used and has achieved great performance recently, where categories are often spatially separated and irrelevant to each other. For example, in ImageNet [71], cats and pens are two distinct categories, sharing few relationships with each other. However, unlike other classification tasks, emotions are strongly correlated and share a unique hierarchical structure with each other according to psychological theories [12], [13].

In our work, we implement Mikel's wheel [12], [13] as base model with eight emotions, *i.e.*, amusement, anger, awe, contentment, disgust, excitement, fear, sad. In addition to its given category, each emotion is assigned with a specific polarity according to Mikel's model [12], [13]. Specifically, amusement, awe, contentment, excitement are positive emotions, while anger, disgust, fear, sad belong to the negative side. Based on the above facts, we design a novel hierarchical cross-entropy loss for VEA, aiming to exploit more emotional-specific information from psychological models. To be specific, we add an auxiliary polarity loss to the traditional cross-entropy loss, with increased penalties on hard false predictions compared with easy ones.

We first send the concatenated emotion feature  $\mathbf{v}_{emo}$  to a classifier and a softmax function successively:

$$p_{emo}(i|\mathbf{v}_{emo}, \mathbf{W}) = \frac{\exp(\mathbf{w}_i \mathbf{v}_{emo})}{\sum_{i=1}^C \exp(\mathbf{w}_i \mathbf{v}_{emo})}, \quad (13)$$

where  $p_{emo} \in \mathbb{R}^C$  represents emotion vector and  $C$  denotes the number of emotion categories in our datasets.  $\mathbf{W} \in \mathbb{R}^{(d_1+d_2+d_3) \times C}$  is a learnable weight matrix of the emotion classifier, which consists of  $\mathbf{w}_i$ .



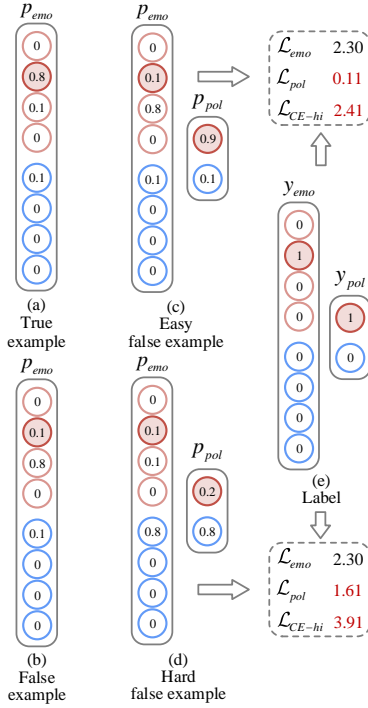


Fig. 7. Numerical examples to validate the effectiveness of the proposed hierarchical cross-entropy loss. It is suggested that emotion loss alone cannot distinguish hard false example from the easy one while the proposed loss is able to separate them apart with an increased penalty.

Since  $p_{emo}$  is a probability vector, we can easily obtain  $\sum_{i=1}^C p_{emo}(i) = 1$ , where each position in  $p_{emo}(i)$  corresponds to the predicted probability of a specific emotion  $i$  in Mikel's model, as shown in Fig. 6. In Fig. 6 (a), take an affective image as an example, we represent the hierarchical labels, *i.e.*, emotion label and polarity label, under the inherent structure of Mikel's wheel. To be specific, the positive emotions (*i.e.*, excitement, amusement, contentment, awe) are correspond to the former  $C/2$  positions in  $p_{emo}$ , while the negative (*i.e.*, sad, fear, disgust, anger) the rest  $C/2$  respectively. Based on the above partition, our emotion vector  $p_{emo}$  can be correspondingly summed up to a polar vector  $p_{pol}$  in accordance with Mikel's model:

$$positive : p_{pol}(0) = \sum_{i=1}^{C/2} p_{emo}(i), \quad (14)$$

$$negative : p_{pol}(1) = \sum_{i=C/2+1}^C p_{emo}(i), \quad (15)$$

where polar vector is also a probability vector that satisfies  $\sum_{j=0}^1 p_{pol}(j) = 1$ . Each value in  $p_{pol}$  can also be viewed as the prediction probability of a specific polarity, *i.e.*, positive, negative. Besides the formulas in Eq. 14–15, we also illustrate the partition process in Fig. 6 (b) for better comprehension.

In traditional cross-entropy loss, there are only two kinds of predictions: true examples (*i.e.*, Fig. 7 (a)) and false ones (*i.e.*, Fig. 7 (b)). However, considering the unique hierarchical structure in emotion classification, false predictions can be divided into simple false examples (*i.e.*, Fig. 7 (c)) and hard false ones

(*i.e.*, Fig. 7 (d)). As shown in Fig. 7, we define simple false examples with emotions incorrectly classified and polarities correctly classified. Similarly, hard false refers to examples whose both emotion and polarity are incorrectly classified. In order to distinguish hard false examples from easy ones, we proposed an auxiliary polarity loss with increased penalties on those hard ones, which is shown in Eq. 16. Meanwhile, true examples and false ones are separated by implementing traditional cross-entropy loss on emotion label, *i.e.*, emotion loss, as shown in Eq. 17:

$$\mathcal{L}_{pol} = \sum_{j=0}^1 y_{pol}(j) \log(p_{pol}(j)), \quad (16)$$

$$\mathcal{L}_{emo} = \sum_{i=1}^C y_{emo}(i) \log(p_{emo}(i)), \quad (17)$$

where  $y_{pol}$  in Eq. 16 represents the polarity label and  $y_{emo}$  in Eq. 17 denotes the emotion label in our datasets.

To illustrate the necessity and effectiveness of the proposed hierarchical cross-entropy loss, we cite a simple example in Fig. 7. Compare Fig. 7 (c) with (e), we can see that with emotion loss alone, the easy false example is treated as identical as the hard one, *i.e.*, with a same numerical result at 2.30. However, polarity loss of two examples are distinguishable from each other, *i.e.*, with a huge discrepancy between 0.11 and 1.61. Thus, the auxiliary polarity loss helps the network to distinguish easy false examples and hard ones apart with increased penalties on hard ones, which is a fine-grained constraint towards precise emotion prediction. Both losses, *i.e.*, emotion loss and polarity loss, are essential and indispensable to the final emotion prediction, which are further combined to form the hierarchical cross-entropy loss:

$$\mathcal{L}_{CE-hi} = \mathcal{L}_{emo} + \lambda \mathcal{L}_{pol}, \quad (18)$$

where  $\lambda$  is a hyper-parameter balancing the importance between the two losses and is further ablated in Sec. IV-D. Finally, the whole network is optimized through Eq. 18 jointly. Therefore, the proposed hierarchical cross-entropy loss can effectively guide the network to pay more attention to hard false examples and consequently helps the network to predict emotions in a hierarchical and emotion-specific manner.

#### IV. EXPERIMENTAL RESULTS

In this section, we first evaluate our approach on four widely-used datasets, including FI [36], EmotionROI [72], Artphoto [20], IAPSA [12], [73], compared with the state-of-the-art methods. Moreover, ablation study is conducted to further validate the effectiveness of our network.

##### A. Datasets

**FI dataset.** The FI dataset [36] is currently the largest well-labeled dataset, containing approximately 23,000 images, which are selected from the Flickr and Instagram by searching eight emotion categories as keywords, *i.e.*, amusement, anger, awe, contentment, disgust, excitement, fear and sad. The weakly labeled emotional images are then well-labeled



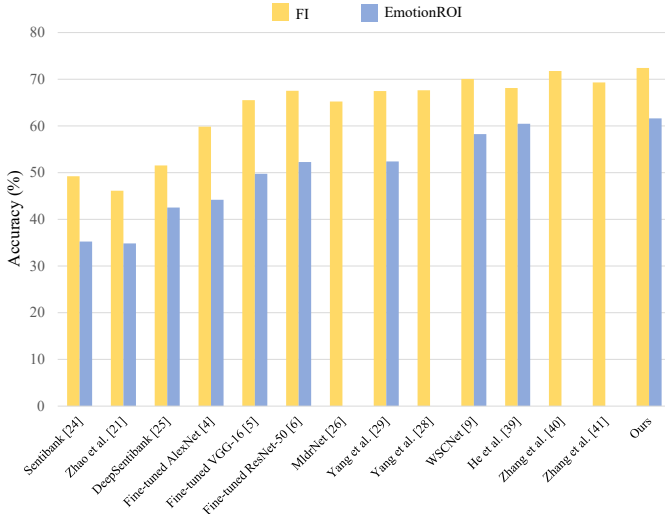


Fig. 8. Classification results on FI dataset and EmotionROI dataset, compared with the state-of-the-art methods.

by 225 Amazon Mechanical Turk (AMT) workers. At least receiving three agreements among five workers, the weakly labels and their corresponding images are preserved as the final FI dataset.

**EmotionROI.** The EmotionROI dataset [72] is a sentiment prediction benchmark collected from Flickr, which contains 1980 images with six emotion categories, *i.e.*, anger, disgust, fear, joy, sad, surprise. Besides categorical label, each image is also annotated on pixel-level with 15 responses corresponding to different emotional regions.

**ArtPhoto.** Images in ArtPhoto dataset [20] are taken by professional artists, aiming to evoke a certain emotion in their photos. There are in total 806 images in this dataset and each of them is labeled with a specific emotion from the eight emotion categories.

**IAPSa dataset.** The International Affective Picture System (IAPS) dataset [73] is a normative emotional stimuli dataset for visual sentiment analysis research. The IAPSa dataset [12] is a subset of the IAPS dataset, which contains 395 images and is labeled with the Mikel's eight sentiment categories.

The above four datasets can be categorized as large-scale dataset (*i.e.*, FI dataset) and small-scale ones (*i.e.*, EmotionROI, ArtPhoto, IAPSa dataset).

### B. Implementation Details

To extract specific emotion features, we initialize our Global-Net with pre-trained parameters from ImageNet [71] while Expression-Net from FER2013 [57]. Besides, our object stimuli detector is initialized with pre-trained parameters from Visual Genome dataset [64], aiming to obtain an accurate description on objects. The whole network is then jointly fine-tuned on emotion datasets with the proposed hierarchical cross-entropy loss in an end-to-end manner.

The FI dataset is randomly split into training set (80%), validation set (5%) and testing set (15%), following the same setting in [36]. The small-scale datasets are split into training set (80%) and testing set (20%) randomly except the one

TABLE I  
CLASSIFICATION ACCURACY (%) ON FI DATASET AND EMOTIONROI DATASET, COMPARED WITH THE STATE-OF-THE-ART METHODS.

Method	FI	EmotionROI
Sentibank [24]	49.23	35.24
Zhao <i>et al.</i> [21]	46.13	34.84
DeepSentibank [25]	51.54	42.53
Fine-tuned AlexNet [4]	59.85	44.19
Fine-tuned VGG-16 [5]	65.52	49.75
Fine-tuned ResNet-50 [6]	67.53	52.27
MldrNet [26]	65.23	–
Yang <i>et al.</i> [29]	67.48	52.40
Yang <i>et al.</i> [28]	67.64	–
WSCNet [9]	70.07	58.25
He <i>et al.</i> [39]	68.13	60.47
Zhang <i>et al.</i> [40]	71.77	–
Zhang <i>et al.</i> [41]	69.32	–
Ours	<b>72.42</b>	<b>61.62</b>

TABLE II  
ABLATION STUDY OF NETWORK STRUCTURE ON FI DATASET.

G-Net		S-Net		E-Net	Acc (%)
RGB	Y	LSTMs	FC layers		
	✓				64.16
✓					67.93
		✓			66.47
				✓	31.64
✓		✓			70.93
✓				✓	70.81
	✓	✓		✓	68.06
		✓		✓	69.69
✓			✓	✓	71.39
✓		✓		✓	<b>72.42</b>

with specified training/testing split, *i.e.*, EmotionROI [72], following the previous work [9], [29], [39]. In addition, each image is resized with its shorter side to 480 and then cropped to 448×448 randomly from the original image or its horizontal flip [6]. We train all our models by using the adaptive optimizer Adam [74], with a weight decay of 5e-5. The learning rate starts from 5e-5 and is decayed by 0.1 every 5 epochs, and the total epoch number is set to 50. Our framework is implemented using PyTorch [75] and our experiments are performed on an NVIDIA GTX 1080ti GPU.

### C. Comparison with the State-of-the-art Methods

1) *Comparison on large-scale dataset:* To evaluate the effectiveness of the proposed method, we first compare our method with the state-of-the-art methods on large-scale dataset FI [36], as shown in TABLE I. Sentibank [24] and Zhao *et al.* [21] adopted hand-crafted features, which were early attempts in VEA. Besides, we also conduct experiments on those widely used CNN backbones, including AlexNet [4], VGG16 [5] and ResNet50 [6], which are initialized with the pre-trained parameters on ImageNet [71] and then fine-tuned on FI. It is obvious that deep learning methods performed

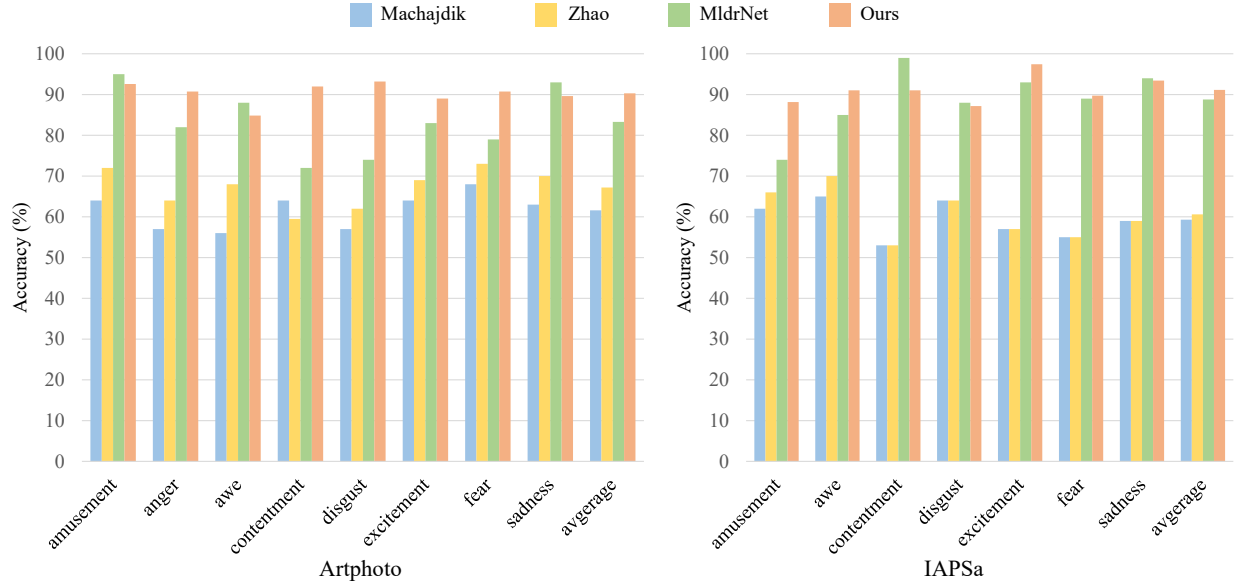


Fig. 9. Classification results on two small-scale datasets, *i.e.*, Artphoto and IAPSA dataset, compared with the state-of-the-art methods.

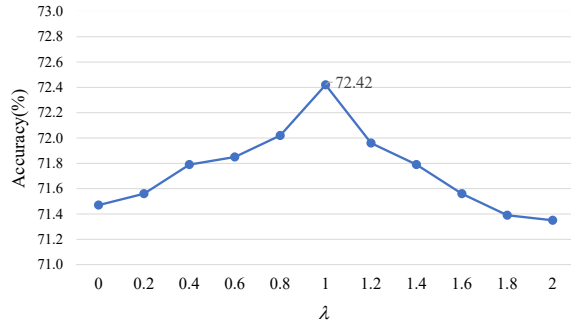


Fig. 10. Hyper-parameter analysis of  $\lambda$  in the hierarchical cross-entropy loss.

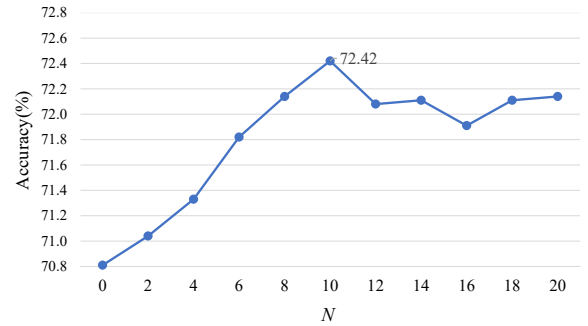


Fig. 11. Hyper-parameter analysis of  $N$  in the stimuli selection process.

favorably against the traditional ones, which attributed to the powerful representation ability of deep features. Based on CNN backbones, researchers subsequently developed specific deep networks to deal with visual emotion analysis, including MldrNet [26], WSCNet [9], Zhang *et al.* [40], Zhang *et al.* [41], *etc.* We further show the holistic curve of classification accuracy in Fig. 8 to summarize the overall performance of the proposed method. Above all, the proposed method outperforms the state-of-the-art methods on FI dataset.

2) *Comparison on small-scale datasets:* We further verify the effectiveness of the proposed method on three small-scale datasets. Considering the class-imbalanced and limited images in IAPSA [73] and Artphoto [20] dataset, we employ the “one against all” strategy to train the classifiers for fair comparison following the previous methods [20]. Besides, image samples in each category are randomly split into five batches and a 5-fold cross-validation is further implemented for the classification results. Since the emotion category of anger only contains eight samples, we remove this category for its sample insufficiency following [20], [21], [26]. As shown in Fig. 9, our method outperforms the state-of-the-art methods, including Machajdik [20], Zhao *et al.* [21] and MldrNet [26].

In addition, we also implement our methods on EmotionROI dataset and validate its superiority towards the state-of-the-art methods, as shown in TABLE I and Fig. 8. It is worth mentioning that the missing data in TABLE I and Fig. 8. is caused by lacking of both classification results and open source codes. Above all, the proposed method consistently outperforms the state-of-the-art methods on both large-scale dataset and small ones, which further proves the effectiveness and robustness of our method.

#### D. Ablation Study

1) *Network architecture analysis:* In TABLE II, we conduct a set of ablation studies to verify the necessity of each stimulus as well as the effectiveness of each module in the proposed network. Our network is divided into three sub-networks, namely Global-Net (G-Net), Semantic-Net (S-Net) and Expression-Net (E-Net). TABLE II shows that when three sub-network acts alone, G-Net reaches the best result, which attributes to the great representation ability of global features. Since not every image in FI dataset contains human faces, the poor performance of E-Net alone is explainable and thus we choose it as an auxiliary branch in the whole network.

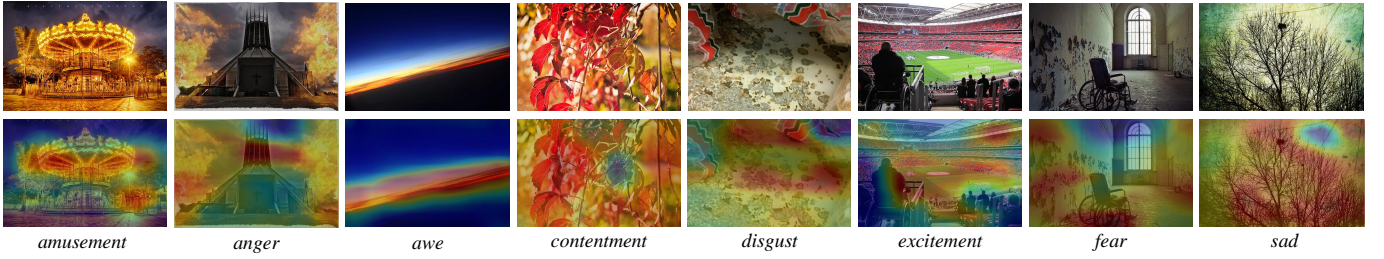


Fig. 12. Visualization of Global-Net. It is shown that Global-Net tend to pay more attention to color features.

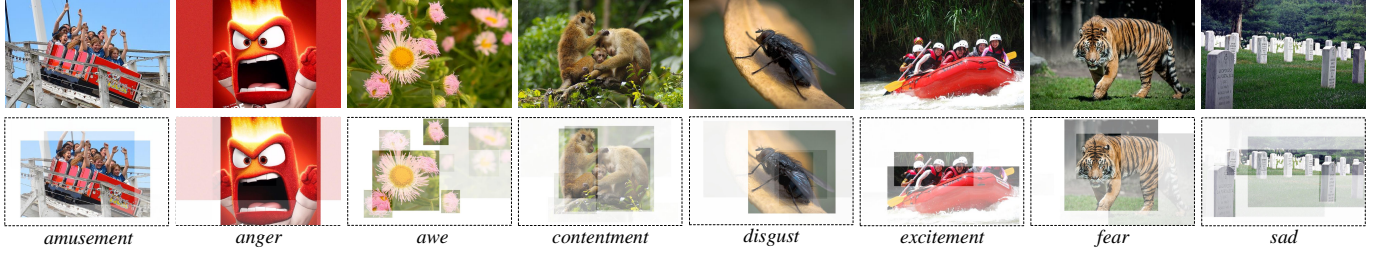


Fig. 13. Visualization of Semantic-Net. Semantic-Net extracts the correlations between different objects, which contributes to the final emotion.

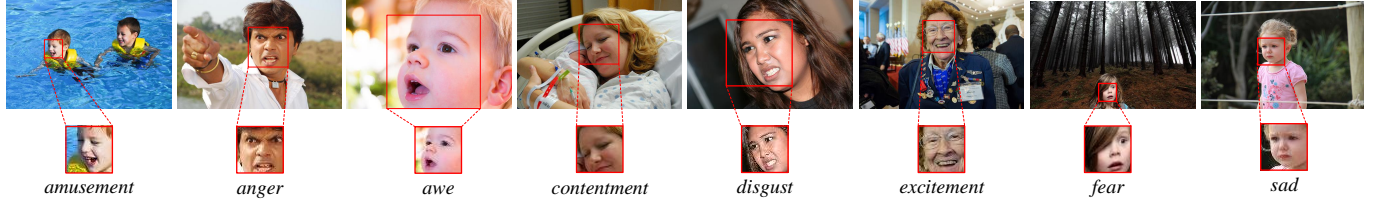


Fig. 14. Visualization of Expression-Net. Expressions in affective images greatly impact the viewer's emotions, which is called empathy.

To validate the necessity of color stimulus, we remove the color information from an image by introducing YCbCr space and restoring Y channel alone. It is obvious that Y channel alone (*i.e.*, 64.16%) has a great performance drop comparing to RGB channels (*i.e.*, 67.93%), which suggests that color features are essential to emotion prediction and G-Net is capable to extract such information. Notably, without color features, the three-branch network (*i.e.*, 69.69%) only outperforms the two-branch network (*i.e.*, 68.06%) to a small extent, where both differ from the final result (*i.e.*, 72.42%) to a large extent. To exploit semantic correlations between different objects, we design two structures for S-Net, namely LSTMs and FC layers, where experimental results show that LSTMs is more capable of extracting semantic correlations. In Sec. III-B2, we demonstrate the superiority of the LSTMs structure theoretically and in this section we prove its effectiveness empirically. We can conclude from TABLE II that each network structure is indispensable and complementary, which jointly contributes to the final result.

2) *Hyper-parameters analysis*: Parameter  $\lambda$  controls the proportion of the two elements in the loss function Eq. (18), which is a decisive hyper-parameter for the classification results. As shown in Fig. 10, the classification accuracy increases as  $\lambda$  grows from 0 to 1 and decreases as  $\lambda$  drops from 1 to 2, in which the peak accuracy reaches 72.42% when  $\lambda = 1$ . When  $\lambda$  is set to 0, the hierarchical cross-entropy loss degrades to

a general cross-entropy loss, which fails to depict the unique hierarchical structure of emotion categories and leads to a sub-optimal result. As  $\lambda$  grows larger, polarity loss dominates the emotion loss, rendering to a decrease in accuracy as well. Hence, we fix  $\lambda = 1$  in the following experimental settings.

The number of bounding boxes, *i.e.*,  $N$ , is another important hyper-parameters in our method. Similar to its common setting, we vary  $N$  from 0 to 20, as shown in Fig. 11. As  $N$  grows from 0 to 10, the classification accuracy is constantly growing. In Sec. III-B2, we assumed that a specific emotion is not evoked by a single objects, but is jointly determined by multiple objects and the correlations between them, which can be further proved by this experiment. However, there exists a slightly drop after 10, due to information redundancy of too many detection boxes. From the above analysis, we choose  $N = 10$  as the node number in our stimuli selection process.

## V. VISUALIZATION RESULTS

To further prove the effectiveness and interpretability of our network, we visualize the three sub-networks (*i.e.*, Global-Net, Semantic-Net, Expression-Net) with examples representing eight different emotions in FI dataset [36]. Besides, we also present some failure cases in our method and analyze the potential reasons behind them.

In Global-Net, we directly extract Class Activation Map (CAM) following [76], where red represents the most attentive





Fig. 15. Failure cases on FI dataset, where red denotes positive emotions and blue represents the negative ones. The former emotion is the labeled one while the latter the incorrectly predicted one.

part while blue the least one. From Fig. 12, we can infer that Global-Net mainly focuses on global features, especially the color feature, when it comes to an image without any distinct objects or faces. Take amusement as an example, Global-Net focuses on the warm lights of carousel, which brings us happiness and joy. The example of fear shows an abandoned hospital, suggesting horror and fear from the color of the whole image.

In Semantic-Net, instead of focusing on a single object, we generate multiple object bounding boxes and leverage attention LSTM and correlation LSTM to weigh the importance and mine the correlations between different objects respectively in Sec. III-B2. In Fig.13, different bounding boxes represent distinct objects, in which the depth of color indicates the importance of the specific object. First, we can infer from Fig.13 that multiple objects jointly determine the final emotion, as opposed to a single one. Besides, different bounding boxes are correlated, where the bounding box with a deeper color suggests an increased attention weight. Fig.13 shows that it is not a single flower but the flower sea amazes us with awe. In the example of excitement, Semantic-Net not only detects happy people, but also the drift, as drifting evokes great excitement indeed.

In Expression-Net, we visualize the cropped faces directly, since facial expressions is the most obvious stimulus to convey emotions, which is called empathy [34]. As mentioned in Sec. III-A, all the faces are cropped to an aligned size of  $48 \times 48$ , shown in the second row of Fig.14. From the visualization results, it is obvious that we are easily affected by the expression given in an image and thus Expression-Net is indispensable in our network. Take contentment as an example, there is a new-born baby lying with his mother, where the satisfied mother may affect us with contentment. Expression-Net detects a crying girl from the example of sad, bringing us sad emotion oppositely.

However, there still exist some failure cases in our method. As shown in Fig. 15, we select some incorrectly classified images from FI dataset, where the former emotion represents the labeled one and the latter the predicted one. Take the first image as an example, the clown boy wears colorful clothes and made up for an exaggerated smile. Considering

these two emotional stimuli, the network tends to predict it with positive emotion, *i.e.*, *excitement*. However, in human cognition, clowns are often disguised to be happy, especially there exists sadness in the boy's eyes, which involves a high-level common sense in human cognition process. Therefore, when encountering complex scenarios, it is hard for deep network to learn such prior knowledge through a simple training process. Besides, some of the incorrect predictions may be attributed to the noisy labels in our datasets, as shown in the below two images in Fig. 15. In FI dataset, most images in *amusement* seem like the former one and most images in *awe* look like the latter one, which makes it understandable to make such false predictions. Above all, failure cases are mainly caused by the complexity and ambiguity of emotions, which may be partly overcome by designing emotion-specific network and labeling datasets by psychologists.

## VI. CONCLUSION

Inspired by S-O-R model in psychology, we propose a stimuli-aware VEA network to simulate the evocation process of human emotions, which consists of stimuli selection, feature extraction and emotion prediction. In the first stage, we introduced stimuli selection into VEA for the first time, by implementing the off-the-shelf tools to choose specific emotional stimuli. After that, we construct three specific sub-networks to extract distinct emotion features from different stimuli simultaneously. Finally, a hierarchical cross-entropy loss is proposed to distinguish hard false examples from easy ones, which helps the network to learn in an emotion-specific manner. Experiments demonstrated that the proposed method consistently outperforms the state-of-the-art methods on four widely-used affective datasets. Ablation study and visualization results further validated the effectiveness and interpretability of the proposed method.

## REFERENCES

- [1] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014, pp. 580–587.
- [2] R. Girshick, "Fast r-cnn," in *ICCV*, 2015, pp. 1440–1448.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015, pp. 91–99.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NeurIPS*, 2012, pp. 1097–1105.
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [7] K. Li, Y. Zhang, K. Li, Y. Li, and Y. Fu, "Visual semantic reasoning for image-text matching," in *ICCV*, 2019, pp. 4654–4662.
- [8] Y. Kao, R. He, and K. Huang, "Deep aesthetic quality assessment with semantic information," *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1482–1495, 2017.
- [9] J. Yang, D. She, Y.-K. Lai, P. L. Rosin, and M.-H. Yang, "Weakly supervised coupled networks for visual sentiment analysis," in *CVPR*, 2018, pp. 7584–7592.
- [10] S. Zhao, G. Ding, Q. Huang, T.-S. Chua, B. W. Schuller, and K. Keutzer, "Affective image content analysis: A comprehensive survey," in *IJCAI*, 2018, pp. 5534–5541.
- [11] R. De Charms, *Personal causation: The internal affective determinants of behavior*. Routledge, 2013.



- [12] J. A. Mikels, B. L. Fredrickson, G. R. Larkin, C. M. Lindberg, S. J. Maglio, and P. A. Reuter-Lorenz, "Emotional category data on images from the international affective picture system," *Behavior research methods*, vol. 37, no. 4, pp. 626–630, 2005.
- [13] S. Zhao, H. Yao, Y. Gao, R. Ji, W. Xie, X. Jiang, and T.-S. Chua, "Predicting personalized emotion perceptions of social images," in *ACMMM*, 2016, pp. 1385–1394.
- [14] H. H. Binali, C. Wu, and V. Potdar, "A new significant area: Emotion detection in e-learning using opinion mining techniques," in *2009 3rd IEEE international conference on digital ecosystems and technologies*, 2009, pp. 259–264.
- [15] A. Yadollahi, A. G. Shahraki, and O. R. Zaiane, "Current state of text sentiment analysis from opinion to emotion mining," *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–33, 2017.
- [16] P. Sánchez-Núñez, M. J. Cobo, C. De Las Heras-Pedrosa, J. I. Peláez, and E. Herrera-Viedma, "Opinion mining, sentiment analysis and emotion understanding in advertising: A bibliometric analysis," *IEEE Access*, vol. 8, pp. 134 563–134 576, 2020.
- [17] C. Chang, "The impacts of emotion elicited by print political advertising on candidate evaluation," *Media psychology*, vol. 3, no. 2, pp. 91–118, 2001.
- [18] A. L. Chapman, "Borderline personality disorder and emotion dysregulation," *Development and Psychopathology*, vol. 31, no. 3, pp. 1143–1156, 2019.
- [19] S. Yang, P. Zhou, K. Duan, M. S. Hossain, and M. F. Alhamid, "emhealth: towards emotion health through depression prediction and intelligent health recommender system," *Mobile Networks and Applications*, vol. 23, no. 2, pp. 216–226, 2018.
- [20] J. Machajdik and A. Hanbury, "Affective image classification using features inspired by psychology and art theory," in *ACMMM*, 2010, pp. 83–92.
- [21] S. Zhao, Y. Gao, X. Jiang, H. Yao, T.-S. Chua, and X. Sun, "Exploring principles-of-art features for image emotion recognition," in *ACMMM*, 2014, pp. 47–56.
- [22] S. Siersdorfer, E. Minack, F. Deng, and J. Hare, "Analyzing and predicting sentiment of images on the social web," in *ACMMM*, 2010, pp. 715–718.
- [23] S. Zhao, H. Yao, Y. Yang, and Y. Zhang, "Affective image retrieval via multi-graph learning," in *ACMMM*, 2014, pp. 1025–1028.
- [24] D. Borth, R. Ji, T. Chen, T. Breuel, and S.-F. Chang, "Large-scale visual sentiment ontology and detectors using adjective noun pairs," in *ACMMM*, 2013, pp. 223–232.
- [25] T. Chen, D. Borth, T. Darrell, and S.-F. Chang, "Deepsentibank: Visual sentiment concept classification with deep convolutional neural networks," *arXiv preprint arXiv:1410.8586*, 2014.
- [26] T. Rao, X. Li, and M. Xu, "Learning multi-level deep representations for image emotion classification," *Neural Processing Letters*, pp. 1–19, 2016.
- [27] J. Yang, D. She, M. Sun, M.-M. Cheng, P. L. Rosin, and L. Wang, "Visual sentiment prediction based on automatic discovery of affective regions," *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2513–2525, 2018.
- [28] J. Yang, D. She, Y.-K. Lai, and M.-H. Yang, "Retrieving and classifying affective images via deep metric learning," in *AAAI*, 2018.
- [29] J. Yang, D. She, and M. Sun, "Joint image emotion classification and distribution learning via deep convolutional neural network," in *IJCAI*, 2017, pp. 3266–3272.
- [30] Q. You, J. Luo, H. Jin, and J. Yang, "Robust image sentiment analysis using progressively trained and domain transferred deep networks," in *AAAI*, 2015.
- [31] N. H. Frijda, *The emotions*. Cambridge University Press, 1986.
- [32] A. Mehrabian and J. A. Russell, *An approach to environmental psychology*. the MIT Press, 1974.
- [33] T. Brosch, G. Pourtois, and D. Sander, "The perception and categorisation of emotional stimuli: A review," *Cognition and emotion*, vol. 24, no. 3, pp. 377–400, 2010.
- [34] C. R. Rogers, "Empathic: An unappreciated way of being," *The counseling psychologist*, vol. 5, no. 2, pp. 2–10, 1975.
- [35] P. Ekman, "An argument for basic emotions," *Cognition & emotion*, vol. 6, no. 3–4, pp. 169–200, 1992.
- [36] Q. You, J. Luo, H. Jin, and J. Yang, "Building a large scale dataset for image emotion recognition: The fine print and the benchmark," in *AAAI*, 2016.
- [37] X. Zhu, L. Li, W. Zhang, T. Rao, M. Xu, Q. Huang, and D. Xu, "Dependency exploitation: A unified cnn-rnn approach for visual emotion recognition," in *IJCAI*, 2017, pp. 3595–3601.
- [38] T. Rao, X. Li, H. Zhang, and M. Xu, "Multi-level region-based convolutional neural network for image emotion classification," *Neurocomputing*, vol. 333, pp. 429–439, 2019.
- [39] X. He, H. Zhang, N. Li, L. Feng, and F. Zheng, "A multi-attentive pyramidal model for visual sentiment analysis," in *IJCNN*, 2019, pp. 1–8.
- [40] W. Zhang, X. He, and W. Lu, "Exploring discriminative representations for image emotion recognition with cnns," *IEEE Transactions on Multimedia*, vol. 22, no. 2, pp. 515–523, 2019.
- [41] J. Zhang, M. Chen, H. Sun, D. Li, and Z. Wang, "Object semantics sentiment correlation analysis enhanced image sentiment classification," *Knowledge-Based Systems*, vol. 191, p. 105245, 2020.
- [42] L. Zhou, Y. Zhang, Y.-G. Jiang, T. Zhang, and W. Fan, "Re-caption: Saliency-enhanced image captioning through two-phase learning," *IEEE Transactions on Image Processing*, vol. 29, pp. 694–709, 2019.
- [43] X. Xiang, M. Dao, G. D. Hager, and T. D. Tran, "Hierarchical sparse and collaborative low-rank representation for emotion recognition," in *ICASSP*, 2015, pp. 3811–3815.
- [44] F. Wang, X. Xiang, C. Liu, T. D. Tran, A. Reiter, G. D. Hager, H. Quon, J. Cheng, and A. L. Yuille, "Regularizing face verification nets for pain intensity regression," in *ICIP*, 2017.
- [45] X. Xiang and T. D. Tran, "Linear disentangled representation learning for facial actions," *IEEE Transactions on Circuits and System for Video Technology*, vol. 28, no. 12, 2018.
- [46] S. Li and W. Deng, "Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 356–370, 2018.
- [47] S. Shojaeilangari, W.-Y. Yau, K. Nandakumar, J. Li, and E. K. Teoh, "Robust representation and recognition of facial emotions using extreme sparse learning," *IEEE Transactions on Image Processing*, vol. 24, no. 7, pp. 2140–2152, 2015.
- [48] P. Chiranjeevi, V. Gopalakrishnan, and P. Moogi, "Neutral face classification using personalized appearance models for fast and robust emotion detection," *IEEE Transactions on Image Processing*, vol. 24, no. 9, pp. 2701–2711, 2015.
- [49] S. Li and W. Deng, "Deep facial expression recognition: A survey," *arXiv preprint arXiv:1804.08348*, 2018.
- [50] D. E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, no. Jul, pp. 1755–1758, 2009.
- [51] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [52] P. Moore, F. Khafa, and L. Barolli, "Semantic valence modeling: emotion recognition and affective states in context-aware systems," in *AINA Workshops*, 2014, pp. 536–541.
- [53] J. Lee, S. Kim, S. Kiim, and K. Sohn, "Spatio-temporal attention based deep neural networks for emotion recognition," in *ICASSP*, 2018, pp. 1513–1517.
- [54] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression," in *CVPR Workshops*, 2010, pp. 94–101.
- [55] M. Pantic, M. Valstar, R. Rademaker, and L. Maat, "Web-based database for facial expression analysis," in *ICME*, 2005, pp. 5–pp.
- [56] M. Valstar and M. Pantic, "Induced disgust, happiness and surprise: an addition to the mmi facial expression database," in *Proc. 3rd Intern. Workshop on EMOTION (satellite of LREC): Corpora for Research on Emotion and Affect*, 2010, p. 65.
- [57] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville, M. Mirza, B. Hamner, W. Cukierski, Y. Tang, D. Thaler, D.-H. Lee *et al.*, "Challenges in representation learning: A report on three machine learning contests," in *ICONIP*, 2013, pp. 117–124.
- [58] C. Fabian Benitez-Quiroz, R. Srinivasan, and A. M. Martinez, "Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild," in *CVPR*, 2016, pp. 5562–5570.
- [59] A. Mollahosseini, B. Hasani, and M. 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, 2017.
- [60] R. S. Woodworth, "Dynamic psychology." 1930.
- [61] O. Buxbaum, "The sor-model," in *Key Insights into Basic Mechanisms of Mental Activity*. Springer, 2016, pp. 7–9.
- [62] X. Cao and J. Sun, "Exploring the effect of overload on the discontinuous intention of social media users: An sor perspective," *Computers in human behavior*, vol. 81, pp. 10–18, 2018.

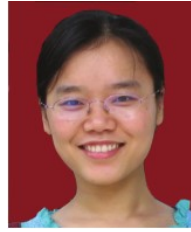
- [63] A. Luqman, X. Cao, A. Ali, A. Masood, and L. Yu, "Empirical investigation of facebook discontinues usage intentions based on sor paradigm," *Computers in Human Behavior*, vol. 70, pp. 544–555, 2017.
- [64] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International Journal of Computer Vision*, vol. 123, no. 1, pp. 32–73, 2017.
- [65] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge 2007 (voc2007) results," 2007.
- [66] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014, pp. 740–755.
- [67] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *CVPR*, 2018, pp. 6077–6086.
- [68] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [69] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015, pp. 2625–2634.
- [70] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [71] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.
- [72] K.-C. Peng, A. Sadovnik, A. Gallagher, and T. Chen, "Where do emotions come from? predicting the emotion stimuli map," in *ICIP*, 2016, pp. 614–618.
- [73] P. J. Lang, M. M. Bradley, B. N. Cuthbert *et al.*, "International affective picture system (iaps): Instruction manual and affective ratings," *The center for research in psychophysiology, University of Florida*, 1999.
- [74] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [75] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [76] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *CVPR*, 2016, pp. 2921–2929.



**Jingyuan Yang** received the B.Eng. degree in Electronic and Information Engineering from Xidian University, Xi'an, China, in 2017. She is currently a Ph. D. Candidate at the School of Electronic Engineering, Xidian University. Her current research interest is visual emotion analysis in deep learning and its applications.



**Jie Li** received the B.Sc. degree in electronic engineering, the M.Sc. degree in signal and information processing, and the Ph.D. degree in circuit and systems from Xidian University, Xi'an, China, in 1995, 1998, and 2004, respectively. She is currently a Professor with the School of Electronic Engineering, Xidian University. Her research interests include image processing and machine learning. In these areas, she has published around 50 technical articles in refereed journals and proceedings, including the IEEE T-IP, T-CSVT, and Information Sciences.



**Xiumei Wang** received the Ph.D. degree from Xidian University, Xi'an, China, in 2010. She is currently a Lecturer with the School of Electronic Engineering, Xidian University. Her current research interests include nonparametric statistical models and machine learning. She has published several scientific articles, including the IEEE Trans. Cybernetics, Pattern Recognition, and Neurocomputing in the above areas.



**Yuxuan Ding** was born in 1995. He received the B.Eng. degree in Intelligent Science and Technology from Xidian University, Xi'an, China, in 2018. He is currently a Ph. D. Candidate at the School of Electronic Engineering, Xidian University. His main research interest covers Machine Learning, Computer Vision, Vision-Language, and their applications.



**Xinbo Gao** (Senior Member, IEEE) received the B.Eng., M.Sc. and Ph.D. degrees in electronic engineering, signal and information processing from Xidian University, Xi'an, China, in 1994, 1997, and 1999, respectively. From 1997 to 1998, he was a research fellow at the Department of Computer Science, Shizuoka University, Shizuoka, Japan. From 2000 to 2001, he was a post-doctoral research fellow at the Department of Information Engineering, the Chinese University of Hong Kong, Hong Kong. Since 2001, he has been at the School of Electronic Engineering, Xidian University. He is a Cheung Kong Professor of Ministry of Education of P. R. China, a Professor of Pattern Recognition and Intelligent System of Xidian University. Since 2020, he is also a Professor of Computer Science and Technology of Chongqing University of Posts and Telecommunications. His current research interests include Image processing, computer vision, multimedia analysis, machine learning and pattern recognition. He has published six books and around 300 technical articles in refereed journals and proceedings. Prof. Gao is on the Editorial Boards of several journals, including Signal Processing (Elsevier) and Neurocomputing (Elsevier). He served as the General Chair/Co-Chair, Program Committee Chair/Co-Chair, or PC Member for around 30 major international conferences. He is a Fellow of the Institute of Engineering and Technology and a Fellow of the Chinese Institute of Electronics.

Engineering, Xidian University. He is a Cheung Kong Professor of Ministry of Education of P. R. China, a Professor of Pattern Recognition and Intelligent System of Xidian University. Since 2020, he is also a Professor of Computer Science and Technology of Chongqing University of Posts and Telecommunications. His current research interests include Image processing, computer vision, multimedia analysis, machine learning and pattern recognition. He has published six books and around 300 technical articles in refereed journals and proceedings. Prof. Gao is on the Editorial Boards of several journals, including Signal Processing (Elsevier) and Neurocomputing (Elsevier). He served as the General Chair/Co-Chair, Program Committee Chair/Co-Chair, or PC Member for around 30 major international conferences. He is a Fellow of the Institute of Engineering and Technology and a Fellow of the Chinese Institute of Electronics.