## Neurocomputing 510 (2022) 107-121

Contents lists available at ScienceDirect

Neurocomputing

journal homepage: www.elsevier.com/locate/neucom

# Cross-individual affective detection using EEG signals with audio-visual embedding



Zhen Liang <sup>a,b</sup>, Xihao Zhang <sup>a,b</sup>, Rushuang Zhou <sup>a,b</sup>, Li Zhang <sup>a,b</sup>, Linling Li <sup>a,b</sup>, Gan Huang <sup>a,b</sup>, Zhiguo Zhang <sup>a,b,c,d,e,\*</sup>

<sup>a</sup> School of Biomedical Engineering, Health Science Center, Shenzhen University, Shenzhen, Guangdong 518060, China

<sup>b</sup> Guangdong Provincial Key Laboratory of Biomedical Measurements and Ultrasound Imaging, Shenzhen, Guangdong 518060, China

<sup>c</sup> Institute of Computing and Intelligence, Harbin Institute of Technology, Shenzhen, China

<sup>d</sup> Marshall Laboratory of Biomedical Engineering, Shenzhen, Guangdong 518060, China

<sup>e</sup> Peng Cheng Laboratory, Shenzhen, Guangdong 518055, China

## ARTICLE INFO

Article history: Received 16 January 2022 Revised 14 May 2022 Accepted 4 September 2022 Available online 9 September 2022 Communicated by Zidong Wang

#### Keywords:

Electroencephalography Individual differences Affective detection Audio-visual embedding Deep domain adaptation

## ABSTRACT

Affective computing is an increasing interdisciplinary research field that provides great potential to recognize, understand and express human emotions. Recently, multimodal analysis starts to gain more popularity in affective studies, which could provide a more comprehensive view of emotion dynamics based on the diverse and complementary information from different data modalities. However, the stability and generalizability of current multimodal analysis methods have not been thoroughly developed yet. In this paper, we propose a novel multimodal analysis method (EEG-AVE: EEG with audio-visual embedding) for cross-individual affective detection, where EEG signals are exploited to identify the emotion-related individual preferences and audio-visual information is leveraged to estimate the intrinsic emotions involved in the multimedia content. EEG-AVE is composed of two main modules. For EEG-based individual preferences prediction module, a multi-scale domain adversarial neural network is developed to explore the shared dynamic, informative, and domain-invariant EEG features across individuals. For video-based intrinsic emotions estimation module, a deep audio-visual feature-based hypergraph clustering method is proposed to examine the latent relationship between semantic audio-visual features and emotions. Through an embedding model, both estimated individual preferences and intrinsic emotions are incorporated with shared weights and further contribute to affective detection across individuals. Experiments on two well-known emotional databases indicate that the proposed EEG-AVE model achieves a better performance under a leave-one-individual-out cross-validation individualindependent evaluation protocol. The results demonstrate that EEG-AVE is an effective model with good reliability and generalizability, which has practical significance in the development of multimodal analysis in affective computing.

© 2022 Elsevier B.V. All rights reserved.

## 1. Introduction

Electroencephalography (EEG) provides a natural way to record human brain activities and has been widely used in the affective intelligence studies [1–5]. In recent years, deep neural network learning methods have provided an effective and efficient approach to characterize informative deep features from EEG data and have achieved promising results in EEG-based affective detec-

\* Corresponding author.

tion applications. For example, a novel dynamic graph convolutional neural network (DGCNN) was proposed in [1] to learn the discriminant and hidden EEG characteristics in a non-linear approach for solving the multi-channel EEG based emotion decoding problem. Jirayucharoensak et al. [6] adopted a stack of several autoencoder structures to perform EEG-based emotion decoding and showed the deep learning network outperformed the traditional classification models such as support vector machine (SVM) and naïve Bayes classifiers. The valid, useful and optimal EEG information can be explored in a deep belief network (DBN) structure, which was demonstrated to be beneficial to the decoding performance [7]. Cui et al. [8] proposed an end-to-end regional-asymmetric convolutional neural network (RACNN) to capture the discriminant EEG features covering temporal, regional,



*E-mail addresses:* janezliang@szu.edu.cn (Z. Liang), zhangxihao2019@email.szu. edu.cn (X. Zhang), 2018222087@szu.edu.cn (R. Zhou), lzhang@szu.edu.cn (L. Zhang), lilinling@szu.edu.cn (L. Li), huanggan@szu.edu.cn (G. Huang), zhiguozhang@hit.edu.cn (Z. Zhang).

and asymmetric information. Based on a series of pretrained stateof-the-art CNN architectures, Cimtay and Ekmekcioglu [4] improved the feature extraction performance and classification capability based on raw EEG signals. Liu et al. [9] proposed a three-dimension convolution attention neural network (3DCANN) to learn the dynamic spatio-temporal joint features with the dual attention weight learning strategy and achieved superior model performance. The existing literature has shown deep learning is a powerful tool in EEG processing, which captures the abstract representations and disentangles the semantic gap between EEG signals and emotional states.

However, due to the problem of individual differences, the stability and generalizability of EEG-based affective detection models are of great challenge. Especially, EEG data are very weak signals and easily susceptible to interference from undesired noises, making it different to distinguish individual-specific and meaningful EEG patterns from noise. The key to solving the problem of individual differences is to minimize the discrepancy in feature distributions across individuals. To improve model generalization to the variance of individual characteristics, transfer learning methods have been introduced and a fruitful line of prior studies has been explored [10-12,2]. Based on feature distribution and classifier parameters learning, Zheng and Lu [11] developed two types of subject-to-subject transfer learning approaches and showed a significant increase in emotion recognition accuracy (conventional generic classifier: 56.73%; the proposed model: 76.31%). Lin and Jung [12] proposed a conditional transfer learning framework to boost a positive transfer for each individual, where the individual transferability was evaluated and effective data from other subjects were leveraged. Li et al. [2] developed a multi-source transfer learning method, where two sessions (calibration and subsequent) were involved and the data differences were transformed by the style transfer mapping and integrated classifier. Among various transfer learning strategies, domain adaptation is a popular way to learn common feature representations and make the feature representations invariant across different domains (source and target domains). Ganin et al. [13] proposed an effective domainadversarial neural network (DANN) to align the feature distributions between the source domain and target domain and also maintain the information of the aligned discriminant features which are predictive of the labels of source samples. Instead of the conventional domain adaptation methods that adapted a well-trained model based on a specific domain to another domain, DANN could well learn the shareable features from different domains and maintain the common knowledge about the given task. Inspired by this work, Li et al. [14] proposed a bihemisphere domain adversarial neural network (BiDANN) model for emotion recognition using EEG signals, in which a global and two local domain discriminators worked adversarially with an emotion classifier to improve the model generalizability. Li et al. [3] proposed a domain adaptation method through simultaneously adapting marginal and conditional distributions based on the latent representations and demonstrated an improvement of the model generalizability across subjects and sessions.

On the other hand, with the great development and application of the internet and multimedia nowadays, there are many approaches to characterize audio-visual content and embed the conveying information with other feature modalities for emotion detection[15–18]. For example, based on traditional handcrafted audio and visual features, Wang et al. [15] investigated several kernel-based methods to analyze and fuse audio-visual features for bimodal emotion recognition. Mo et al. [16] proposed Hilbert-Huang Transform (HHT) based visual and audio features for a time–frequency-energy description of videos and introduced cross-correlation features to indicate the dependencies between the visual and audio signals. Furthermore, the recent success of deep learning methods in computer vision brings new insights into the video-content-based affective study. Acar et al. [19] utilized CNNs to learn mid-level audio-visual feature representations for affective analysis of music video clips. Zhang et al. [17] proposed a hybrid deep model to characterize a joint audio-visual feature representation for emotion recognition, where CNN, 3D-CNN, and DBN were integrated with a two-stage learning strategy.

In general, current affective computing models can be mainly categorized into two streams. One stream is to predict individual preferences through analyzing a user's spontaneous physiological responses (i.e. EEG signals) while watching the videos [20,21]. The individualized reactions to emotions are well-considered, and an assumption is made here that different emotions could be elicited for different viewers when watching the same video. However, spontaneous response-based individual preferences prediction would be sensitive to individual differences and fail to achieve reliable performance in affective detection across individuals. Another stream is to estimate intrinsic emotions from video content itself by integrating visual and audio features in either feature-level fusion or decision fusion and building a classifier for distinguishing emotions [17,22]. The video content-based intrinsic emotions estimation could achieve a stable emotion detection performance, but it fails to consider the deviations of individuals in emotion perception. This motivates us to study the underlying associations among emotions, video content, and brain responses, where video content functions as a stimulation clue indicating what kind of emotions would possibly be elicited and brain responses reveal individual emotion perceiving processes showing how we exactly feel the emotions. An appropriate embedding strategy of individual preferences and intrinsic emotions in cross-individual affection detection tasks could be helpful to learn reliable affective features from video content and benefit in enhancing the estimation stability of individual emotions.

Besides, compared to unimodal analysis, multimodal fusion could provide more details, compensate for the incomplete information from another modality, and develop advanced intelligent affective systems [23]. Recently, Wang et al. [24] incorporated video information and EEG signals to improve the video emotion tagging performance. This study characterized a set of traditional visual and audio features, including brightness, color energy, and visual excitement for visual features, and average energy, average loudness, spectrum flow, zero-crossing rate (ZCR), the standard deviation of ZCR, 13 Mel-Frequency Cepstral Coefficients (MFCC) and the corresponding standard deviations for audio features. The proposed hybrid emotion tagging approach was realized on a modified SVM classifier, and the corresponding performance was improved from 54.80% to 75.20% for valence and from 65.10% to 85.00% for arousal after a fusion of multi-modality data. Inspired by the success of the embedding protocol across different data modalities, this study proposes a novel affective information detection model (termed EEG-AVE) to learn transferable features from EEG signals individual preferences prediction with an embedding of affective-related multimedia characteristics (intrinsic emotions estimation) to enhance cross-individual affective detection performance. The proposed EEG-AVE model is illustrated in Fig. 1, which is composed of three parts: EEG-based individual preferences prediction, audio-visual based intrinsic emotions estimation, and multimodal embedding.

• **EEG-based individual preferences prediction.** In this part, we propose a multi-scale domain adversarial neural network (termed as MsDANN hereinafter) based on DANN [13] to enhance the generalization ability of EEG feature representation across individuals and boost the model performance on individual preferences prediction. Specifically, EEG data from different individuals are treated as domains, where the source domain



Fig. 1. The proposed EEG-AVE model.

refers to the existing individuals and the target domain refers to the new coming individual(s). Based on the input multi-scale feature representation, the feature extractor network, task classification network, and discriminator network are designed to make the source and target domains share similar and close latent distribution to work with the same prediction model. As the mining of emotional informative and sensitive features from EEG signals is still a great challenge, this study introduces a multi-scale feature representation to improve feature efficacy and model adaptability to complex and dynamic emotion cases. Compared to single-scale feature representation, pioneer studies have shown EEG signals analysis with such a coarse-grain procedure could be beneficial to emotion studies [25–27].

- Audio-visual based intrinsic emotions estimation. To enhance the model stability in cross-individual affective detection tasks, audio-visual content analysis is conducted to digest the intrinsic emotion information involved in the videos which could be used as supplementary information for individual affective detection. Due to the well-known "semantic gap" or "emotional gap" that the traditional handcrafted features may fail to sufficiently discriminate emotions, we develop a deep audio-visual feature-based hypergraph clustering method (termed as DAVFHC) for characterizing semantic and high-level audio-visual features. Here, two pretrained CNN architectures (VGGNet [28] and VGGish [29], whose performance have been widely recognized in audio-visual information analysis [30,31]) are adopted to explore the emotion-related audiovisual characteristics and the most optimal features are fused through a hypergraph theory.
- **Multimodal embedding.** The final affective detection result is determined by an embedding model where the predicted individual preferences from EEG signals and the estimated intrinsic emotions from audio-visual content are fused at a decision level. The compensation information from different modalities contributes together to tackling the individual differences problems in affective detection.

The major contributions of this work are summarized as follows. (1) We propose a novel cross-individual affective detection model (EEG-AVE) to incorporate spontaneous brain responses and stimulation clues in a hybrid embedding strategy. Both EEG and audio-visual information are exploited to digest different dimensions of emotions, and the compensation relationships among different data modalities in the affective detection study are examined. (2) We introduce an effective individual preferences prediction method (MsDANN) to estimate individual emotions from EEG signals, where the impact of individual differences is diminished through a transfer learning approach. (3) We present an efficient intrinsic emotions estimation method (DAVFHC) to characterize the emotion-related in audio-visual materials as the supplementary information for the cross-individual affective study. Here, the semantic audio-visual features are extracted by using deep learning methods, and the complex and latent relationships of deep audio-visual features with emotion labels are measured with hypergraph theory.

## 2. Methodology

## 2.1. Individual preferences prediction

In this section, we propose a new transfer learning based neural network, MsDANN, to address the problem of the individual differences in EEG-based emotion detection. In this network, a multiscale feature representation is incorporated to capture a series of rich feature characteristics of EEG signals and maximize the informative context for predicting a diverse set of individual preferences in emotions. Specifically, we extract the differential entropy (DE) features [32] from the defined frequency sub-bands (refer to Table 1) at different frequency/scale resolutions (1 Hz, 0.5 Hz, and 0.25 Hz), and build respective domain adaptation models with domain adversarial training methods. In the proposed MsDANN, the common features from different individuals are learned; at the same time, the relationships between the learned common features and the related emotional information are preserved. The network structure of MsDANN is shown in Fig. 2, which is composed of three parts: the generator (feature extractor network) for deep feature extraction, the classifier (task classification

Z. Liang, X. Zhang, R. Zhou et al.

## Table 1

The defined frequency sub-bands for DE feature characterization.

	$\theta$	α1	α2	β1	β2	β3	γ1	γ2	γ3
frequency band (Hz)	4-8	8-10	10-13	13-16	16-20	20–28	28-34	34-39	39-45



Fig. 2. The proposed MsDANN model.

network) for emotion label prediction, and the discriminator (discriminator network) for real or fake data distinguishing. The corresponding network configurations about the designed MsDANN are reported in Table 2. Here, the generator and classifier could be considered as a standard feed-forward architecture, while the generator and discriminator are trained based on a gradient reversal layer to ensure the feature distributions of two domains as indistinguishable as possible. In this study, the EEG data with emotion labels are treated as the source domain to train the generator, classifier, and discriminator; while the EEG data without emotion labels are utilized to train the generator and discriminator. Through this multi-scale deep framework, a set of transferable features involving affective information could be characterized, the cross-domain discrepancy could be bridged, and the classification

#### Table 2

The network configurations of MsDANN.

	Name	Input Size	Output Size
Generator	Full Connection	32×50	32×4
	ELU Activation Function	32×4	32×4
	Flatten	32×4	128×1
	Full Connection	128×1	64×1
	ELU Activation Function	64×1	64×1
	Full Connection	64×1	64×1
	ELU Activation Function	64×1	64×1
Classifier	Full Connection	64×1	2×1
	Dropout	2×1	2×1
	Softmax Activation Function	2×1	2×1
Discriminator	Full Connection	64×1	128×1
	<b>ReLU Activation Function</b>	128×1	128×1
	Full Connection	128×1	1×1
	Sigmoid Activation Function	$1 \times 1$	1×1

performance could be effectively improved in both source and target domains.

To learn a shared common feature space between the source and target domains and also guarantee the learned feature representation involving enough information for revealing the emotion states, the loss objective function is designed below. Suppose that the source and target domains are denoted as  $\mathbb{S}$  and  $\mathbb{T}$ . In the domain learning, the EEG data with emotion labels in S are given as  $x^{l} = \{x_{1}^{l}, \dots, x_{N_{S}}^{l}\}$  and  $y = \{y_{1}, \dots, y_{N_{S}}\}$ , where  $x_{i}^{l}$  is the input EEG data at *l*th scale feature representation and  $y_i$  is the corresponding emotion label of  $x_i^l$ ,  $N_s$  is the sample size of  $x^l$ . On the other hand, the unlabeled EEG data in  $\ensuremath{\mathbb{T}}$  is denoted as  $z^{l} = \{z_{1}^{l}, \dots, z_{N_{T}}^{l}\}$ , where  $z_{i}^{l}$  is the input EEG data at *l*th scale feature representation and  $N_T$  is the corresponding sample size of  $z^l$ . We denote the generator, classifier, and discriminator as  $r_{\theta}$ ,  $c_{\sigma}$ ,  $d_{\mu}$  with the parameters of  $\theta, \sigma$  and  $\mu$ . To ensure the learned features by  $r_{\theta}$ from the source domain or target domain are indistinguishable, the domain adversarial training objective function is given as

$$\min_{\theta} \max_{\mu} \mathbf{E}_{(x^{l}, z^{l}) \sim (\mathbb{S}, \mathbb{T})} \mathscr{L}_{D}(\mu, \theta, x^{l}) + \mathscr{L}_{D}(\mu, \theta, z^{l}),$$
(1)

where  $\mathscr{L}_D$  is a binary cross-entropy loss for the discriminator to be trained to distinguish  $\mathbb{S}$  and  $\mathbb{T}$ , defined as

$$\begin{aligned} \mathscr{L}_{D}(\mu,\theta,\mathbf{x}^{l}) &= -\mathbb{I}[\mathbf{x}^{l}\sim\mathbb{S}]\log(d_{\mu}\circ r_{\theta}(\mathbf{x}^{l})) - \mathbb{I}[\mathbf{x}^{l}\sim\mathbb{T}] \\ &\times \log(1-d_{\mu}\circ r_{\theta}(\mathbf{x}^{l})). \end{aligned}$$
(2)

Here, I is an indicator function, and  $\circ$  refers to the composite mapping from the generator to the discriminator. Based on Eq. (1), we add another loss function  $\mathscr{L}_T$  for the classifier part as

$$\min_{\boldsymbol{\sigma},\boldsymbol{\theta}} \max_{\boldsymbol{\mu}} \mathbf{E}_{\mathbf{x}^{l} \sim \mathbb{S}} \left[ \mathscr{L}_{T}(\boldsymbol{\sigma},\boldsymbol{\theta},\mathbf{x}^{l}) \right] + \lambda \mathbf{E}_{\left(\mathbf{x}^{l} \mathbf{z}^{l}\right) \sim (\mathbb{S},\mathbb{T})} \left[ \mathscr{L}_{D}(\boldsymbol{\mu},\boldsymbol{\theta},\mathbf{x}^{l}) + \mathscr{L}_{D}(\boldsymbol{\mu},\boldsymbol{\theta},\mathbf{z}^{l}) \right],$$
(3)

where  $\mathscr{L}_T(\sigma, \theta, x^l)$  is the classification loss in the source domain, determined by  $\sum Loss(c_{\sigma} \circ r_{\theta}(x^l), y)$ .  $\lambda$  is a balance parameter during the learning process, given as

$$\lambda = \frac{2}{1 - \exp(-\gamma p)} - 1, \tag{4}$$

where  $\gamma$  is a constant value and *p* is a factor of an epoch. Eq. (3) is the final objective function for MsDANN model training. The proposed MsDANN model is an end-to-end framework for crossindividual emotion prediction based on EEG signals, combining the feature learning adaptation and emotion classification into a unified deep model. Based on the input data with multi-scale DE feature representation, the domain adaptation and classification loss are exploited to guide the generator to learn effective feature representations across individuals via the gradient reversal layer and efficiently tackle the problem of the individual differences in EEG data processing.

## 2.2. Intrinsic Emotions Estimation

At present, a number of well trained deep CNN models have been successfully applied to multimedia processing, such as Alex-Net [33], GoogLeNet [34] and VGG [28] for visual content, and VGGish [29] for audio content. The deep features could bridge the semantic gap and improve semantic interpretation performance. In this section, we develop a DAVFHC method to learn and decode the semantic features from audio-visual content for intrinsic emotions estimation.

At the visual level, a pretrained VGGNet network [28] is utilized to process frame-based visual information and characterize effective visual features. The training and testing data sets were based on ILSVRC-2012, with 1.3 M training pictures, 50 K test pictures, and 100 K validation pictures. The network was trained by optimizing a polynomial logistic regression objective function with the smallest batch-based gradient descent momentum. Considering the balance of layer depth and performance, VGG16 is utilized in this paper to characterize the frame-based visual features. It consists of 13 convolutional layers and 3 fully connected layers. The corresponding number of convolution kernels at each layer are 64, 64, 128, 128, 256, 256, 256, 512, 512, 512, 512, 512, 512, and 512, and the kernel size is  $3 \times 3$ . As illustrated in Fig. 3, the visual feature extraction procedure includes three steps.

- 1. **1. Frame-based visual feature extraction.** The video frames are input to the pretrained VGG16 and the corresponding feature maps are characterized at each convolutional layer. For each layer, an average feature map is then calculated and converted into a feature vector.
- 2. 2. Segment-based visual feature extraction. Instead of direct averaging all the frame-based features in one segment, we introduce an adaptive keyframe detection step to detect a keyframe from every segment based on the feature distribution. Suppose that one segment is composed of k frames with the corresponding extracted features, denoted as  $B^{l} = \{b_{1}^{l}, \dots, b_{k}^{l}\},\$ where  $i = 1, ..., N_i$  refers to the convolutional layer. The keyframe detection is illustrated as follows. (1) All frames are grouped into one cluster in terms of  $B^{i}$ ; (2) The cluster center  $c^{i}$  is computed; (3) The distance between each frame  $b_{i}^{i}$  $(i \in [1, k])$  and the cluster center  $c^i$  is calculated, denoted as  $\{d_1^i, \ldots, d_k^i\}$ ; (4) the frame which is the closest to  $c^i$  is selected as the keyframe of the segment, termed as  $k^* = \arg\min\{d_1^i, \dots, d_k^i\}$ . Then, the corresponding feature of the keyframe  $b_{\nu^*}^l$  is treated as the segment-based feature representation.
- 3. **3. Segment-based visual feature fusion.** The characterized segment-based features at each single convolutional layer  $(b_{k^*}^i, i \in [1, N_i])$  are then fused by concatenation. Empirically, the segment length is set to 1s. To get the semantic features, only the characterized features at the last two convolutional layers (i = 12 and 13) are used as visual features ( $\Psi_V$ ) in the proposed DAVFHC method.

At the audio level, a pretrained CNN network, VGGish [29], is adopted to characterize effective audio features. VGGish is a deep network model trained on a Youtube-8 M database (training/vali-



Fig. 3. The visual feature extraction procedure.

dation/test: 70 M/20 M/10 M), which has been proved to be capable of extracting effective and efficient deep auditory features in various applications [35,36,31]. The network contains 6 convolutional layers, and the corresponding numbers of convolution kernels are 64, 128, 256, 256, 512, and 512, respectively. The kernel size is  $3 \times 3$ . Same as the visual feature extraction process, the audio features are also characterized at the segment level. The audio feature extraction procedure totally includes four steps below.

- 1. **Data preparation.** The audio signals are detected from the emotional clips and then partitioned into a number of segments with a fixed length.
- 2. **Data preprocessing.** The segment-based audio data is preprocessed following the procedures presented in [29].
- 3. **Deep audio feature characterization.** For each segment, the logarithmic Mel spectrum is characterized and input to the VGGish. The deep feature maps are extracted at each convolutional layer and averaged into one feature map.
- 4. **Deep audio feature fusion.** For each segment, the feature map at every single convolutional layer is converted into a feature vector. The converted feature vectors across different convolutional layers are then fused by concatenation. Empirically, the segment length is set to 1s (same as the visual data. To get the semantic features, only the feature vectors extracted from the last two layers (5th and 6th) are used as audio features ( $\Psi_A$ ) in the proposed DAVFHC method.

The characterized segment-based visual and audio features are concatenated and formed into a segment-based audio-visual feature vector termed  $\Psi_M = [\Psi_V, \Psi_A]$ . The complex relationships among all the segments from the emotional clips are constructed with a hypergraph which has been widely recognized as an effective approach for complex hidden data structure description. For the traditional graph, only pairwise relationships between any two vertices are considered, which would lead to the information loss [37]. In the hypergraph, one edge (termed as hyperedge in the hypergraph) could connect more than two vertices and the complex relationship among a group of vertices could be well described. In the paper, the segments are the vertices denoted as *V*, and the connections among the segments are the hyperedges denoted as *E*. One hypergraph could be represented as G = (V, E), where the vertices and hyperedges are denoted as  $V = \{v_1, v_2, \dots, v_{|V|}\}$  and  $E = \{e_1, e_2, \dots, e_{|E|}\}$ , respectively. |V| and |E| are the corresponding vertex size and hyperedge size. Here, the vertices belong to one hyperedge  $e_k \in E$  is termed as  $\left\{ v_1^{e_k}, v_2^{e_k}, \dots, v_{|e_k|}^{e_k} \right\}$ . To define the vertices and hyperedges relationships, the similarity between any two vertices ( $v_i$  and  $v_i$ ) are measured in terms of the segment-based audio-visual feature representation  $(\Psi_M^{\nu_i} = \left\{\psi_{M,1}^{\nu_i}, \dots, \psi_{M,N_M}^{\nu_i}\right\}$  and  $\Psi_M^{\nu_j} =$  $\left\{\psi_{M,1}^{v_j},\ldots,\psi_{M,N_M}^{v_j}
ight\}$ ), as

$$a\left(\Psi_{M}^{\nu_{i}},\Psi_{M}^{\nu_{j}}\right) = \frac{1}{1 + \xi_{\Psi_{M}^{\nu_{i}},\Psi_{M}^{\nu_{j}}}},\tag{5}$$

where  $N_M$  is the feature dimensionality, and  $\xi_{\Psi_M^{\nu_i},\Psi_M^{\nu_j}}$  is the calculated distance, given as

$$\xi_{\Psi_{M}^{\nu_{i}},\Psi_{M}^{\nu_{j}}} = \sum_{t=1,\dots,N_{M}} \frac{\left(\psi_{M,t}^{\nu_{i}} - \psi_{M,t}^{\nu_{j}}\right)^{2}}{\psi_{M,t}^{\nu_{i}} + \psi_{M,t}^{\nu_{j}}}.$$
(6)

Based on the measured similarity matrix  $A = \left\{ a \left( \Psi_M^{v_i}, \Psi_M^{v_j} \right) \right\}_{i,j=1}^N$ (*N* is the sample size), an incidence matrix  $H = \left\{ h(v_i, e_k) \right\}_{i,k}^{|V|,|E|}$  is formed, in which the connection relationships between the vertices V and the hyperedges E is described as

$$h(v_i, e_k) = \begin{cases} 1 & \text{if } v_i \in e_k \\ 0 & \text{if } v_i \notin e_k \end{cases}.$$
(7)

The hyperedge weight matrix,  $W = diag(w(e_1), \ldots, w(e_k), \ldots, w(e_{|E|}))$ , is a diagonal matrix indicating the weights of all the hyperedges *E* in the hypergraph *G*. The weight  $w(e_k)$  of one hyperedge  $e_k \in E$  is computed based on the calculated similarities among the vertices that belong to  $e_k$ , given as

$$w(e_k) = \frac{\sum_{v_i, v_j \in e_k, v_i \neq v_j} a\left(\Psi_M^{v_i}, \Psi_M^{v_j}\right)}{\tau},\tag{8}$$

where  $a(\Psi_M^{\nu_i}, \Psi_M^{\nu_j})$  is the similarity value between the vertices of  $\nu_i$ and  $v_i$ , given in Eq. (5).  $\tau$  is the total number of vertices connected to the hyperedge  $e_k$ . As  $w(e_k)$  is a measurement of all the similarity relationships among the vertices that belong to one hyperedge, a higher  $w(e_k)$  value indicates a strong connection of homogeneous vertices of the hyperedge and a lower  $w(e_k)$  refers to a weak connection of the hyperedge in which the connected vertices share little similar properties. In other words, the hypergraph structure could well describe the relationships of the audio-visual segments in terms of properties. The vertex degree matrix.  $D_{\nu} = diag(d(\nu_1), \dots, d(\nu_i), \dots, d(\nu_{|V|}))$ , is a diagonal matrix presenting the degree of all the vertices in the hypergraph G. The degree of one vertex  $v_i \in V$  is calculated as the summation of all the hyperedge weights of the hyperedges (e) that the vertex belongs to, defined as

$$d(v_i) = \sum_{e \in E \mid v_i \in e} h(v_i, e) w(e).$$
(9)

The hyperedge degree matrix,  $D_e = diag(d(e_1), \ldots, d(e_k), \ldots, d(e_{|E|}))$ , is also a diagonal matrix showing the degree of all the hyperedges in the hypergraph *G*. The degree of one hyperedge  $e_k \in E$  is calculated as the summation of all the vertices (v) that connect to the hyperedge, given as

$$d(e_k) = \sum_{e_k \in E | v \in e_k} h(v, e_k).$$
(10)

In this study, we introduce a spectral hypergraph partitioning method [38] to partition the constructed hypergraph into a number of clusters corresponding to the emotion states (high or low). Thus, it is a two-way hypergraph partitioning problem that could be described as

$$Hcut\left(S,\overline{S}\right) = \sum_{e \in \partial S} w(e) \frac{|e \cap S||e \cap \overline{S}|}{d(e)},\tag{11}$$

where *S* and  $\overline{S}$  are the partitions of the vertices *V*. For two-way partitioning,  $\overline{S}$  is the complement of *S*.  $\partial S$  is the partition boundary, given as  $\partial S = \left\{ e \in E | e \cap S \neq \emptyset \text{ and} e \cap \overline{S} \neq \emptyset \right\}$ . w(e) is the hyperedge weight calculated in Eq. (8), and d(e) is the hyperedge degree defined in Eq. (10). To avoid unbalanced partitioning,  $Hcut\left(S,\overline{S}\right)$  is further normalized by

$$NHcut\left(S,\overline{S}\right) = Hcut\left(S,\overline{S}\right)\left(\frac{1}{vol(S)} + \frac{1}{vol(\overline{S})}\right),\tag{12}$$

where vol(S) and  $vol(\overline{S})$  are the volumes of S and  $\overline{S}$ , given as  $vol(S) = \sum_{v \in S} d(v)$  and  $vol(\overline{S}) = \sum_{v \in \overline{S}} d(v)$ . Here, d(v) is the vertex degree given in Eq. (9). The partitioning rule is to look for the weak-

est hyperedge *e* between *S* and  $\overline{S}$ , where the vertices in the same cluster should be tightly connected (high hyperedge weights) and the vertices in the different clusters should be weakly connected (low hyperedge weights). An optimal partitioning is given in Eq. (13) to find the weakest connection between two partitions, which is an NP-complete problem solved by a real-valued optimization method. *f* is a label vector to be learned, which contains the affective clustering information.

$$\arg \min_{f} \frac{1}{2} \sum_{e \in E} \sum_{v_{i}, v_{j} \in V} \frac{w(e)h(v_{i}, e)h(v_{j}, e)}{d(e)} \\ \times \left( \frac{f(v_{i})}{\sqrt{d(v_{i})}} - \frac{f(v_{j})}{\sqrt{d(v_{j})}} \right)^{2} \\ = \arg \min_{f} \sum_{e \in E} \sum_{v_{i}, v_{j} \in V} \frac{w(e)h(v_{i}, e)h(v_{j}, e)}{d(e)} \\ \times \left( \frac{f^{2}(v_{i})}{\sqrt{d(v_{i})}} - \frac{f(v_{i})f(v_{j})}{\sqrt{d(v_{i})d(v_{j})}} \right) \\ = \arg \min_{f} \sum_{v_{i} \in V} f^{2}(v_{i}) \sum_{e \in E} \frac{w(e)h(v_{i}, e)}{d(v_{i})} \sum_{v_{j} \in V} \frac{h(v_{j}, e)}{d(e)} \\ - \sum_{e \in E} \sum_{v_{i}, v_{j} \in V} \frac{f(v_{i})h(v_{i}, e)w(e)h(v_{j}, e)f(v_{j})}{\sqrt{d(v_{i})d(v_{j})d(e)}}$$
(13)

$$= \arg \min_{f} f^{T} (I - \Theta) f$$

where  $\Theta$  is given as

$$\Theta = D_{v}^{-(1/2)} HW D_{e}^{-1} H^{T} D_{v}^{-(1/2)}, \qquad (14)$$

and *I* is an identity matrix with the same size as *W*. Here,  $D_v$ , H, W,  $D_e$  are the vertex degree matrix (Eq. (9)), the incidence matrix (Eq. (7)). the weight matrix (Eq. (8)), and the hyperedge degree matrix (Eq. (10)) defined above. The hypergraph Laplacian is denoted as

$$\Delta = I - \Theta. \tag{15}$$

The optimal solution is transformed to find the eigenvectors of  $\Delta$  whose eigenvalues are the smallest. In other words, the optimal hypergraph partitioning results find the top eigenvectors with the smallest non-zeros eigenvalues in  $\Delta$  and form an eigenspace for the subsequent vertex clustering with the K-means method. Through this approach, all the vertices are grouped into two clusters. The corresponding emotional state of each cluster is determined by the majority distribution of the involved vertices. If most vertices belong to the high level, the cluster's emotion state is assigned as high; on the other hand, it is assigned as low. In practice, to avoid information leaking, the clusters' emotional states are only determined based on the training samples.

## 2.3. Embedding Model

Based on the aforementioned work, we incorporate the estimated intrinsic emotions based on deep audio-visual features and the predicted individual preferences from the collected simultaneous EEG signals, and conduct a decision-level information fusion for final affective prediction. Specifically, we fuse EEG signals and audio-visual information in a decision level through shared weights. Suppose that the predicted emotional individual preferences based on EEG signals are denoted as  $Y^{EEG} = \{y_1^{EEG}, \ldots, y_N^{EEG}\}$  and the estimated intrinsic emotions based on audio-visual content are denoted as  $Y^{Video} = \{y_1^{Video}, \ldots, y_N^{Video}\}$ . The final detected affective results are determined by

$$y_i^{FUS} = \frac{w^{EEG} \times y_i^{EEG} + w^{Video} \times y_i^{Video}}{w^{EEG} + w^{Video}},$$
(16)

where  $w^{EEG}$  and  $w^{Video}$  are the shared weights of EEG signals and audio-visual information in the fusion process.  $Y^{FUS} = \{y_1^{FUS}, \ldots, y_N^{FUS}\}$  are the final affective detection results. In the implementation, considering the individual preferences  $(y_i^{EEG})$ predicted from EEG signals and the intrinsic emotions  $(y_i^{Video})$  estimated by multimedia information play the same important role in the cross-individual affective detection, the shared weights  $(w^{EEG}$  and  $w^{Video})$  are set equally.

## 3. Experimental Results

In this section, we conduct extensive experiments on MAHNOB-HCI [39] and DEAP [40] databases which are commonly used to evaluate the effectiveness of cross-individual affective studies. To crosscompare with other studies, two types of groundtruth data which are commonly used in the literature are adopted here to evaluate the experimental results. One is the aggregated groundtruth, where different participants watching one video are tagged with the same emotion label. Another is the **non-aggregated ground**truth, where different participants watching one video are tagged with different emotion labels according to the corresponding subjective assessment. Different from the aggregated groundtruth, different participants would have different emotional feelings about the same video, due to the differences in background, experience, religion, education, and so on. In other words, the non-aggregated groundtruth could be more capable of reflecting the emotional dynamics in individuals and should be more encouraged to be used for affective detection evaluation. During the training process, the stochastic gradient descent (SGD) is adopted, with the training epoch, learning rate, momentum, minibatch size, and weight decay of 1000, 0.001, 0.9, 48, and 0.001, respectively. To fully evaluate the model generalizability, a strict leave-one-individual-out crossvalidation is introduced through the model evaluation process. The model is implemented on an NVIDIA GeForce RTX 2080 GPU, with CUDA 10.0 using the Pytorch API. The source code is available at https://github.com/KAZABANA/EEG-AVE.

#### 3.1. Emotional EEG Databases

The MAHNOB-HCI database [39] contains EEG data of 30 participants (male/female: 13/17; age:  $26.06\pm4.39$ ) from different cultural backgrounds. A total of 20 commercial film clips (duration: from 34.9s to 117s, with an average of 81.4s and a standard deviation of 22.5s) were selected for emotional eliciting. After the emotional clip played, the participants were requested to give a subjective assessment of their emotions during watching the emotional clip using a score in the range of 1 to 9. During the experiment, EEG signals were simultaneously collected at a sampling rate of 256 Hz, by using the Biosemi active II system with 32 Ag/AgCl electrodes placed according to the standard international 10–20 electrode system. Due to the data incompleteness of participants 3, 9, 12, 15, 16, and 26, only 24 participants are used in this paper.

The DEAP database [40] consists of 32 subjects' EEG emotion data. A total of 40 music videos, with a fixed duration of the 60s, were selected for emotional eliciting. The corresponding subjective feedback on different emotion dimensions was collected for each music video. The EEG signals were recorded at a sampling rate of 512 Hz from 32 active AgCl electrode sites according to the international 10–20 system placement.

## 3.2. Experiment Protocols

To cross-compare with the results presented in the other studies, we utilize a fixed threshold of 5 for scores (in the range of 1 to 9) to discretize the subjective feedback into high and low levels ( $\geq$  5 high; < 5 low) as the non-aggregated groudtruth. The aggregated groundtruth is an average of all the returned subjective feedback for one video. Two performance metrics, detection accuracy  $P_{acc}$  and F1-Score  $P_f$ , are used to validate the evaluation performance.  $P_{acc}$  is an overall detection performance measurement and  $P_f$  is a harmonic average of the precision and sensitivity which is less susceptible to the unbalanced classification problems. The corresponding definitions are given as

$$P_{acc} = \frac{n_{TN} + n_{TP}}{n_{TN} + n_{FN} + n_{TP} + n_{FP}} \times 100\%, \tag{17}$$

and

$$P_f = \frac{2 \times P_{pre} \times P_{sen}}{P_{pre} + P_{sen}} \times 100\%, \tag{18}$$

where  $n_{TN}$  and  $n_{TP}$  are the correctly predicted samples, and  $n_{FN}$  and  $n_{FP}$  are the incorrectly predicted samples. The precision  $P_{pre}$  and sensitivity  $P_{sen}$  are given as

$$P_{pre} = \frac{n_{TP}}{n_{TP} + n_{FP}},\tag{19}$$

$$P_{\text{sen}} = \frac{n_{TP}}{n_{TP} + n_{FN}}.$$
(20)

To fully evaluate the validity and reliability of the model performance, a strict leave-one-out cross-validation is adopted. All the predicted individual preferences and the estimated intrinsic emotions are obtained in a cross-validation manner. For the proposed MsDANN model, the model training and testing are conducted on a leave-one-individual-out cross-validation. In one round of cross-validation, all the samples from 1 individual are treated as the test data, while the other samples from the remaining individuals are used as the training data. Until each participant is treated as the test data once, the final result of MsDANN is a formation of all the obtained test results through the cross-validation rounds. For the developed DAVFHC method, the model training and testing are conducted on a strict leave-one-video-out cross-validation. In one round of cross-validation, all the samples from 1 video are used as test data and the other samples from the remaining videos are treated as training data. Until each video is treated as the test data once, the final prediction result of DAVFHC is a formation of the obtained test results in all the cross-validation rounds. In other words, after obtaining all the test results of all EEG and video samples in the above-mentioned cross-validation rounds, the final affective results are obtained by a decision fusion.

#### 3.3. Cross-Individual Affective Detection Experiments

To improve the affective detection performance, both EEG signals and audio-visual information are embedded in the proposed EEG-AVE model. Here, we roughly estimate what kind of emotion could be triggered according to the audio-visual content itself (intrinsic emotions estimation), and detect the individual preferences of each individual through analyzing the recording EEG signals while he/she is watching the multimedia material (individual preferences detection). The contributions of EEG signals and audiovisual information through the affective detection process are considered equally important. The corresponding emotion decoding performance for valence and arousal on MAHNOB-HCI and DEAP databases are reported in Table 3. We compare EEG-AVE model with the existing representative methods such as [39,41–43,24]. It is worth noting that the experimental results presented in [24] were evaluated with the aggregated groundtruth.

For the MAHNOB-HCI database, our proposed model outperforms the existing methods for valence, where the  $P_{acc}$  and  $P_{f}$ results are 90.21% and 90.45% for the aggregated groundtruth and 71.13% and 66.83% for non-aggregated groundtruth. For the results with non-aggregated groundtruth, even the obtained  $P_{acc}$ values of our proposed EEG-AVE model and Rayatdoost and Soleymani [43]'s work are comparable, a better  $P_f$  of our proposed EEG-AVE model is observed, where  $P_f$  is 62.08% for Rayatdoost and Soleymani [43]'s work and 66.83% for our model (improved by 7.65%). For the results with aggregated groundtruth, our proposed EEG-AVE model increases the affective detection performance by 19.96% for  $P_{acc}$  and 22.56% for  $P_f$ , compared to Wang et al. [24]'s work (Pacc: 75.20%; Pf: 73.80%). Similar promising emotion recognition performance is observed for arousal. The  $P_{acc}$  and  $P_f$  results of our proposed model are 85.59% and 86.55% for the aggregated groundtruth and 66.47% and 63.25% for non-aggregated groundtruth. For aggregated groundtruth, the proposed EEG-AVE model performs better than Wang et al. [24] (*P*<sub>acc</sub>: 85.00%; *P*<sub>f</sub>: 82.40%), with the corresponding increase rate of 0.69% ( $P_{acc}$ ) and 5.04%  $(P_f)$ . For non-aggregated groundtruth, the EEG-AVE model also gains better performance than the existing methods on recognition accuracy. Besides, the above results show aggregated groundtruth leads to a higher detection performance compared to the nonaggregated groundtruth, as the individual differences in emotional feelings about the clip are not considered.

For the DEAP database, our proposed model outperforms the existing methods for valence in terms of both accuracy and F1-score. For aggregated groundtruth, the  $P_{acc}$  and  $P_f$  results are 75.26% and 77.16%; For non-aggregated groundtruth, the  $P_{acc}$  and  $P_f$  results are 68.50% and 68.81%. Compared to Wang et al. [24]'s work ( $P_{acc}$ : 71.10%;  $P_f$ : 68.60%), our model increases the affective detection performance on valence with the increase rates of 5.85% ( $P_{acc}$ ) and 12.48% ( $P_f$ ). Even the affective detection accuracy for arousal is not as good as the existing methods, where  $P_{acc}$  values are 71.92% and 54.52% for aggregated groundtruth and non-aggregated groundtruth, respectively. The obtained F1-score values are the highest, where  $P_f$  values are 79.50% and 60.80% for aggregated groundtruth and non-aggregated groundtruth and non-aggre

## Table 3

Affective detection performance on MAHNOB-HCI and DEAP databases.

Methods	Groundtruth	MAHNOB-HCI				DEAP			
		Valence		Arousal		Valence		Arousal	
		Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$
Soleymani et al. [39]	Non-Aggregated	57.00	56.00	52.40	42.00	-	-	-	-
Zhu et al. [41]	Non-Aggregated	58.16	56.36	61.35	63.08	-	-	-	-
Huang et al. [42]	Non-Aggregated	62.13	-	61.80	-	-	-	-	-
Rayatdoost and Soleymani [43]	Non-Aggregated	71.25	62.08	61.46	50.60	59.22	56.68	55.70	50.02
Wang et al. [24]	Aggregated	75.20	73.80	85.00	82.40	71.10	68.60	79.00	69.20
Proposed EEG-AVE model	Aggregated	90.21	90.45	85.59	86.55	75.26	77.16	71.92	79.50
Proposed EEG-AVE model	Non-Aggregated	71.13	66.83	66.47	63.25	68.50	68.81	54.52	60.80

[44], F1-score is a better and more important metric for classification models which can distinguish specific types of errors including false positives and false negatives.

## 4. Discussion and Conclusion

To fully study the EEG-AVE performance, we also compare the proposed model with different embedding strategies and domain adaptation conditions. Besides, we also examine the effect of deep and handcrafted multimedia affective representations.

## 4.1. Performance Evaluation of Embedding Strategy

We compare the affective detection performance when different embedding strategies are adopted. Here are three embedding strategies: EEG + Visual + Audio (the proposed EEG-AVE model), EEG + Visual (only visual information embedded with EEG signals), and EEG + Audio (only audio information embedded with EEG signals). The corresponding affective detection performances for valence and arousal with aggregated and non-aggregated groundtruth on MAHNOB-HCI and DEAP databases are summarized in Table 4.

The results of the MAHNOB-HCI database show that EEG-based affective detection with an embedding of both visual and audio information achieves the best performance for both valence and arousal. For EEG + Visual strategy, the affective detection performance for valence decreases to 74.65% (aggregated) and 67.75% (non-aggregated) for  $P_{acc}$  and 74.61% (aggregated) and 61.79% (non-aggregated) for  $P_f$ ; while the affective detection performance for arousal decreases to 77.28% (aggregated) and 63.26% (nonaggregated) for Pacc and 78.57% (aggregated) and 59.27% (nonaggregated) for P<sub>f</sub>. The average decrease rates of valence and arousal are 11.76% and 7.51%, respectively. For EEG + Audio embedding strategy, the affective detection performance for valence decreases from 90.21% to 69.08% for  $P_{acc}$  and from 90.45% to 73.06% for  $P_f$ when aggregated groundtruth is utilized; while it decreases from 71.13% to 58.57% for  $P_{acc}$  and from 66.83% to 58.27% for  $P_f$  when non-aggregated groundtruth is used. A similar decrease pattern is also observed on the affective detection performance for arousal, where it decreases from 85.59% to 68.55% for  $P_{acc}$  and from 86.55% to 72.20% for  $P_f$  when aggregated groundtruth is adopted; while it decreases from 66.47% to 54.91% for Pacc and from 63.25% to 53.64% for  $P_f$  when non-aggregated groundtruth is utilized. The average decrease rates of valence and arousal are 18.28% and 17.27%, respectively. The comparison results with different embedding strategies reveal that an embedding of both visual and audio information has a capability to reach better affective detection performance, compared to only visual or audio embedded. In addition, we find only visual embedded outperforms only audio embedded, which suggests that visual information plays a more critical role in emotion perception, especially in film clips.

For the DEAP database, EEG-based affective detection with an embedding of both visual and audio information achieves the best performance for valence. For aggregated groundtruth, the  $P_{acc}$  and  $P_f$  values decrease from 75.26% and 77.16% (EEG + Visual + Audio) to 62.68% and 65.20% (EEG + Visual) and to 71.46% and 74.58% (EEG + Audio). The average decrease rate is 10.15%. For nonaggregated groundtruth, the  $P_{acc}$  and  $P_f$  values decrease from 68.50% and 68.81% (EEG + Visual + Audio) to 63.32% and 63.29% (EEG + Visual) and to 62.11% and 63.89% (EEG + Audio). The average decrease rate is 8.02%. However, for affective detection of arousal, similar results are obtained for EEG + Audio and EEG + Visual + Audio. One possible reason could be that the embedding strategies of visual and audio information could be different for different emotional dimensions. For example, it is observed that compared to visual information, audio plays a more important role for affective detection in the DEAP database, as the used stimuli for emotion-evoking were music videos. For the MAHNOB-HCI database, the affection detection performance is more relied on visual information, as the used stimuli for emotion eliciting were movie clips. To further validate that incorporating both individual preferences and intrinsic emotions could be beneficial to identify human emotion changes, we also evaluate the pure EEG-based detection performance. For aggregated groundtruth, the affective detection performance is dramatically enhanced. The EEG based affective detection results are 55.21% (valence) and 50.78% (arousal) for MAHNOB-HCI database and 55.27% (valence) and 47.33% (arousal) for DEAP database. While, the EEG + Visual + Audio based affective detection results are 90.21% (valence) and 85.59% (arousal) for MAHNOB-HCI database and 75.26% (valence) and 71.92% (arousal) for DEAP database. For non-aggregated groundtruth, the corresponding affective detection results with pure EEG signals are 69.15% (valence) and 66.93% (arousal) for MAHNOB-HCI database. For DEAP database, the corresponding affective detection results with pure EEG signals are 64.97% (valence) and 66.78% (arousal). Comparing EEG and EEG + Visual + Audio results, the detection performance on valence increases from 69.15% to 71.13% (with an increased rate of 2.86%) for MAHNOB-HCI database and from 64.97% to 66.78% (with an increased rate of 5.43%) for DEAP database.

#### 4.2. Performance Evaluation of Domain Adaptation Effect

To analyze the domain adaptation effect in solving the problem of the individual differences, we also introduce a baseline method, a multi-scale neural network (termed as MsNN shown in Fig. 4, with the network configurations in Table 5) for model comparison under the condition without deep domain adaptation. Based on the input data with multi-scale DE feature representation, no feature adaptation or transfer learning is adopted between the source domain and target domain. In other words, the difference in the feature distribution extracted from different individuals is not considered through modeling. The corresponding loss function is given as

## Table 4

Affective detection performance with different embedding strategies on MAHNOB-HCI and DEAP databases.

	Embedding Strategy	MAHNOB-	HCI			DEAP			
		Aggregated		Non-Aggregated		Aggregated		Non-Aggregated	
		Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$
Valence	EEG + Visual	74.65	74.61	67.75	61.79	62.68	65.20	63.32	63.29
	EEG + Audio	69.08	73.06	58.57	58.27	71.46	74.58	62.11	63.89
	EEG + Visual + Audio	90.21	90.45	71.13	66.83	75.26	77.16	68.50	68.81
Arousal	EEG + Visual	77.28	78.57	63.26	59.27	65.22	78.10	43.22	58.82
	EEG + Audio	68.55	72.20	54.91	53.64	72.37	79.59	55.50	61.10
	EEG + Visual + Audio	85.59	86.55	66.47	63.25	71.92	79.50	54.52	60.80



Fig. 4. The baseline MsNN model for model comparison.

Table 5

The network configurations of MsNN.

	Name	Input Size	Output Size
Generator	Full Connection	32×50	32×4
	ELU Activation Function	32×4	32×4
	Flatten	32×4	128×1
	Full Connection	128×1	64×1
	ELU Activation Function	64×1	64×1
	Full Connection	64×1	64×1
	ELU Activation Function	64×1	64×1
Classifier	Full Connection	64×1	2×1
	Dropout	2×1	2×1
	Softmax Activation Function	2×1	$2 \times 1$

$$\min_{\sigma,\theta} \mathbf{E}_{\mathbf{x}^{l} \sim \mathbb{S}} \left[ \mathscr{L}_{T} \left( \sigma, \theta, \mathbf{x}^{l} \right) \right], \tag{21}$$

where  $\mathbf{E}_{\mathbf{x}' \sim \mathbb{S}} [\mathscr{Q}_T(\sigma, \theta, \mathbf{x}')]$  is the cross-entropy loss function in the source domain (training data). Compared to the loss function of MsDANN given in Eq. (3), the domain adversarial loss function between the source and target domains is removed in Eq. (21). For model validation of MsNN, same as the proposed MsDANN, the EEG-based emotional individual preferences prediction is trained and tested on the source domain and target domain separately, under a strict leave-one-video-out cross-validation protocol. The corresponding affective detection performance of MsDANN and MsNN based EEG-AVE model for valence and arousal detection on MAHNOB-HCI and DEAP databases are reported in Table 6.

For aggregated results on the MAHNOB-HCI database, compared to MsDANN based EEG-AVE model under the embedding strategy of EEG + Visual + Audio, the detection performance of MsNN based EEG-AVE model decreases by 9.10% and 7.19% in terms of  $P_{acc}$  and  $P_f$ , respectively. Comparing MsDANN and MsNN based EEG-AVE model performance under EEG + Visual and EEG + Audio embedding strategies, both  $P_{acc}$  and  $P_f$  values also have similar decrease patterns. The  $P_{acc}$  value decreases from 74.65% to 70.04% for EEG + Visual, and from 69.08% to 65.33% for EEG + Audio. The  $P_f$  value declines from 74.61% to 73.03% for EEG + Visual, and from 73.06% to 71.64% for EEG + Audio. When non-aggregated groundtruth is used, cross-comparing MsDANN and MsNN based model performance under an embedding strategy of EEG + Visual + Audio, it is found that the decoding performance significantly decreases from 71.13% to 63.58% (decreased by

10.61%) in terms of  $P_{acc}$  and from 66.83% to 62.35% (decreased by 6.7%) in terms of  $P_f$ . Similar trends are also observed in the other embedding strategies. The corresponding detection accuracies decrease to 60.98% ( $P_{acc}$ ) and 59.20% ( $P_f$ ) for EEG + Visual embedding strategy, and to 53.38% ( $P_{acc}$ ) and 56.34% ( $P_f$ ) for EEG + Audio embedding strategy. On the other hand, for arousal detection, the aggregated results show MsDANN based EEG-AVE model outperforms MsNN based EEG-AVE model across all three different embedding strategies in terms of both  $P_{acc}$  and  $P_f$ . For EEG + Visual + Audio, EEG + Visual and EEG + Audio embedding strategies, the corresponding improvement rates from MsNN to MsDANN are 9.18%, 7.30% and 4.93% for  $P_{acc}$ , and that are 6.30%, 3.71% and 1.08% for P<sub>f</sub>. For non-aggregated results, similar patterns are observed. Better results are achieved when MsDANN based EEG-AVE model is adopted. Here, the improvement rates for three different embedding strategies are 14.60%, 14.73% and 10.73% for P<sub>acc</sub> and 8.73%, 7.90%, and 2.72% for P<sub>f</sub>. Similar comparison results are also observed on the DEAP database, where MsDANN generally performs better than MsNN across all three embedding strategies (EEG + visual, EEG + Audio, EEG + Visual + Audio). For example, comparing MsDANN and MsNN based EEG-AVE model performance under EEG + Visual + Audio embedding strategy, the  $P_{acc}$ and P<sub>f</sub> values of valence decrease from 75.26% and 77.16% to 66.03% and 72.55% for aggregated groundtruth and from 68.50% and 68.81% to 55.59% and 61.84% for non-aggregated groundtruth. The average decrease rate is 11.80%. For arousal, MsDANN outperformed MsNN when non-aggregated groundtruth is adopted.

Besides, to further verify the domain adaption effect under different feature representations (**single-scale** and **multi-scale** feature representations), we take the MAHNOB-HCI database as an example and evaluate the individual preferences detection results on individuals when the single-scale (low-scale/ middle-scale/ and high-scale) and multi-scale (a fusion of low-scale, middle-scale, and high-scale) feature representations are adopted. Here, the performance comparisons are conducted based on MsDANN and MsNN, and the corresponding comparison results are reported in Table 8 and 9. For both MsDANN- and MsNN-based individual preferences detection models, the multi-scale feature representation (a fusion of low-scale, middle-scale, and high-scale) consistently performs superiorly against the single-scale feature representations (low-scale/ middle-scale/ high-scale) at the individual level and group level. For MsDANN, the cross-individual model performance

#### Table 6

Affective detection performance of MsDANN and MsNN on MAHNOB-HCI and DEAP databases using deep features.

	Embedding Strategy	EEG Model	MAHNOB-HCI				DEAP	_		
			Aggregated		Non-Aggregated		Aggregated		Non-Aggregated	
			Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$
Valence	EEG + Visual	MsDANN	74.65	74.61	67.75	61.79	62.68	65.20	63.32	63.29
		MsNN	70.04	73.03	60.98	59.20	58.21	66.07	52.15	58.69
	EEG + Audio	MsDANN	69.08	73.06	58.57	58.27	71.46	74.58	62.11	63.89
		MsNN	65.33	71.64	53.38	56.34	64.47	71.67	53.02	60.21
	EEG + Visual + Audio	MsDANN	90.21	90.45	71.13	66.83	75.26	77.16	68.50	68.81
		MsNN	82.00	83.95	63.58	62.35	66.03	72.55	55.59	61.84
Arousal	EEG + Visual	MsDANN	77.28	78.57	63.26	59.27	65.22	78.10	43.22	58.82
		MsNN	72.02	75.76	55.14	54.93	65.20	78.11	42.45	58.30
	EEG + Audio	MsDANN	68.55	72.20	54.91	53.64	72.37	79.59	55.50	61.10
		MsNN	65.33	71.43	49.59	52.22	72.27	79.58	49.96	56.42
	EEG + Visual + Audio	MsDANN	85.59	86.55	66.47	63.25	71.92	79.50	54.52	60.80
		MsNN	78.39	81.42	58.00	58.17	72.62	80.07	50.62	57.59

#### Table 7

Affective detection performance of MsDANN and MsNN on MAHNOB-HCI and DEAP databases using handcrafted features.

Embedding Strategy	EEG Model	MAHNOB	-HCI			DEAP				
		Aggregated		Non-Ag	gregated	Aggre	gated	Non-Aggregated		
		Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	
EEG + Visual	MsDANN	62.68	64.75	51.38	46.23	59.96	58.80	61.40	57.04	
	MsNN	59.96	65.08	47.45	47.03	56.51	62.75	50.66	54.90	
EEG + Audio	MsDANN	64.11	62.70	59.51	49.86	59.31	65.92	55.62	60.39	
	MsNN	61.82	64.42	54.41	50.37	56.52	66.89	48.74	58.65	
EEG + Visual + Audio	MsDANN	73.52	69.26	64.99	50.48	61.05	63.65	59.89	59.82	
	MsNN	68.75	68.58	58.43	50.52	57.23	65.08	50.25	56.79	
EEG + Visual	MsDANN	70.04	71.99	55.23	50.87	64.56	76.64	46.73	59.27	
	MsNN	66.08	71.15	50.08	50.91	65.06	77.02	44.69	57.83	
EEG + Audio	MsDANN	65.23	71.41	52.95	55.51	62.00	71.44	53.94	58.91	
	MsNN	62.51	70.97	48.03	54.11	62.99	72.36	48.77	54.64	
EEG + Visual + Audio	MsDANN	72.18	72.92	59.79	53.71	61.69	67.74	58.93	58.00	
	MsNN	67.33	71.39	52.80	51.99	63.11	69.34	50.89	50.58	
	Embedding Strategy EEG + Visual EEG + Audio EEG + Visual + Audio EEG + Visual EEG + Audio EEG + Visual + Audio	Embedding Strategy EEG Model EEG + Visual MsDANN EEG + Audio MsDANN EEG + Visual + Audio MsDANN EEG + Visual + Audio MsDANN EEG + Audio MsDANN EEG + Audio MsDANN EEG + Audio MsDANN EEG + Visual + Audio MsDANN MsNN	Embedding StrategyEEG ModelMAHNOBAggreAggrePaccPaccEEG + VisualMsDANN62.68MSNN59.96EEG + AudioMsDANN64.11MSNN61.82EEG + Visual + AudioMSDANN73.52MSNN68.75MSNN68.75EEG + VisualMSDANN70.04MSNN66.08MSDANN65.23EEG + AudioMSDANN62.51EEG + Visual + AudioMSDANN62.51EEG + Visual + AudioMSDANN72.18MSNN67.33MSNN	$\begin{tabular}{ c c c c } \hline EEG Model & \underline{MAHNOB-HCl} \\ \hline \\ $	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{c c c c c c c } \mbox{EEG Model} & \underline{MAHNOB-HCl} & \underline{Aggregated} & \underline{Non-Aggregated} \\ \hline & Aggregated & P_{f} & P_{acc} & P_{f} \\ \hline & P_{acc} & P_{f} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{f} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{f} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P_{acc} \\ \hline & P_{acc} & P_{acc} & P_{acc} & P$	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{c c c c c c c c c c } EEG Model & \underline{MAHNOB-HCl} & \underline{Pacc} & Non-Aggregated & Aggregated & A$	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	

#### Table 8

MsDANN-based individual preferences detection results of each individual (*S*1, *S*2, ..., *S*24) under leave-one-individual-out cross-validation on MAHNOB-HCI database. Here, **left-top:** low-scale feature representation; **right-top:** middle-scale feature representation; **left-bottom:** high-scale feature representation; **right-bottom:** the proposed multi-scale feature representation. The average performance of low-scale, middle-scale, high-scale, and multi-scale feature representations are 64.44%, 64.51%, 63.96%, and 69.16%, respectively.

S1	S2	S3	S4	S5	S6	S1	S2	S3	S4	S5	S6
65.01	73.74	70.23	59.96	67.40	62.92	65.74	74.54	70.17	59.23	65.38	64.39
S7	S8	S9	S10	S11	S12	S7	S8	S9	S10	S11	S12
56.03	71.83	66.97	70.11	68.51	68.14	56.21	71.03	68.33	70.23	69.86	68.02
S13	S14	S15	S16	S17	S18	S13	S14	S15	S16	S17	S18
69.13	54.24	56.70	52.95	62.30	63.78	66.42	54.43	58.73	52.64	60.89	64.82
S19	S20	S21	S22	S23	S24	S19	S20	S21	S22	S23	S24
63.35	67.16	69.37	58.92	72.63	55.23	64.33	68.51	68.70	59.84	69.31	56.58
S1	S2	<b>S</b> 3	S4	S5	<i>S</i> 6	S1	S2	S3	S4	S5	S6
63.90	68.70	69.00	59.41	65.74	63.35	68.82	77.49	71.83	60.76	70.79	67.10
S7	S8	S9	S10	S11	S12	S7	S8	S9	S10	S11	S12
58.36	68.94	66.61	69.74	68.57	68.51	60.15	76.51	78.35	74.60	72.88	73.06
S13	S14	S15	S16	S17	S18	S13	S14	S15	S16	S17	S18
71.65	53.69	56.21	50.06	57.81	65.87	78.60	56.46	61.07	58.73	64.15	66.73
S19	S20	S21	S22	S23	S24	S19	S20	S21	S22	S23	S24
60.82	66.67	70.17	59.29	74.78	57.26	65.13	72.57	77.49	62.24	85.36	58.92

results (group level) using the single-scale feature representation are 64.44% (low-scale), 64.51% (middle-scale), and 63.96% (highscale), while the corresponding performance based on the multiscale feature representation enhances to 69.16%. The results show that, compared to the single-scale feature representation, the multi-scale feature representation could be beneficial to the affective detection performance, with an average increase rate of 7.55%. A similar phenomenon is observed in MsNN-based individual preference detection results. The single-scale feature representation results are 60.43% (low-scale), 59.92% (middle-scale), 60.50% (high-scale), and the multi-scale feature representation result is 61.41% (the average increase rate is 1.87%). The above results confirm the proposed multi-scale feature representation is robust enough for EEG decoding. On the other hand, comparing the results presented in Table 8 and 9, the benefit of domain adaption is consistently observed under both single-scale and multi-scale feature

#### Table 9

MsNN-based individual preferences detection results of each individual (*S1*, *S2*, ..., *S24*) under leave-one-individual-out cross-validation on MAHNOB-HCI database. Here, **left-top:** low-scale feature representation; **right-top:** middle-scale feature representation; **left-bottom:** high-scale feature representation; **right-bottom:** the proposed multi-scale feature representation. The average performance of low-scale, middle-scale, high-scale, and multi-scale feature representations are 60.43%, 59.92%, 60.50%, and 61.41%, respectively.

S1	S2	S3	<i>S</i> 4	S5	S6	S1	S2	S3	S4	S5	S6
62.42	60.33	58.98	54.98	55.78	63.59	53.63	66.05	54.86	56.46	58.00	66.30
S7	S8	S9	S10	S11	S12	S7	S8	S9	S10	S11	S12
57.13	68.51	49.45	72.45	70.42	56.70	52.58	67.71	59.16	67.40	68.57	57.20
S13	S14	S15	S16	S17	S18	S13	S14	S15	S16	S17	S18
76.57	51.97	54.67	50.74	59.41	52.15	75.58	51.97	56.33	47.60	58.67	49.82
S19	S20	S21	S22	S23	S24	S19	S20	S21	S22	S23	S24
57.20	64.70	71.77	60.82	57.44	62.18	59.72	65.25	68.82	58.12	56.77	61.44
S1	S2	S3	S4	S5	<i>S</i> 6	S1	S2	S3	S4	<i>S5</i>	S6
56.46	64.51	61.50	54.74	58.49	66.73	58.06	65.31	59.96	55.10	58.18	66.73
S7	S8	S9	S10	S11	S12	S7	S8	S9	S10	S11	S12
58.49	65.19	55.78	63.90	68.82	58.12	56.33	68.57	55.97	71.16	71.16	58.61
S13	S14	S15	S16	S17	S18	S13	S14	S15	S16	S17	S18
73.55	51.11	55.60	58.36	53.75	55.84	76.69	52.28	56.03	53.38	59.10	52.34
S19	S20	S21	S22	S23	S24	S19	S20	S21	S22	S23	S24
55.23	64.45	70.30	56.52	63.78	60.82	58.06	66.11	72.32	59.35	59.66	63.28

representations for each individual. These results further demonstrate the reliability and generalizability of the proposed MsDANN. Note that all the results are measured under non-aggregated groundtruth.

The above results demonstrate that, compared to MsNN, MsDANN is much more powerful in the proposed EEG-AVE model to deal with the problem of the individual differences in EEG signal processing. It provides a reliable and useful way to adaptively learn the shared emotion-related common and discriminant feature representation across individuals and demonstrates the validity of the domain adaptation method in EEG-based affective detection applications.

## 4.3. Performance Evaluation of Multimedia Representation

In this study, audio-visual information is represented by deep features characterized by two pretrained networks. We further verify the effectiveness of the deep feature representation and compare it with the performance using more traditional handcrafted features. Inspired from the previous video affective studies [45–47,24], the commonly used handcrafted features are extracted and compared here. For visual information representation, the adopted handcrafted features include lighting key features, color information, and shadow portions in the HSL and HSV spaces. For audio information representation, the used traditional audio features include energy, loudness, spectrum flux, zero-crossing rate (ZCR), Mel-frequency cepstral coefficients (MFCCs), log energy, and the standard deviations of the above ZCR, MFCC, and log energy. The affective analysis of multimedia content with different feature representations is conducted and the corresponding comparison results of valence and arousal are summarized in Table 7. The results show that, compared to the performance presented in Table 6, a significant improvement in affective detection performance is obtained when deep feature representation is used instead of handcrafted features. It reveals that compared to the traditional handcrafted feature representation, deep feature representation is a better and richer affective representation for understanding and perceiving the multimedia content.

## 4.4. Performance Evaluation of Fusion Manner

As shown in Eq. (16), the final affective detection results are determined by a fusion of EEG information and video content, referring to the individual preferences and intrinsic emotions, respectively. This fusion manner could be termed **EEG-Video**-

**Fusion**. To further verify the efficiency and effectiveness of the presented EEG-Video-Fusion, we also evaluate the affective detection performance under another fusion manner, by separately incorporating EEG information, visual content, and audio content as

$$y_{i}^{FUS} = \frac{w^{EEG} \times y_{i}^{EEG} + w^{Visual} \times y_{i}^{Visual} + w^{Audio} \times y_{i}^{Audio}}{w^{EEG} + w^{Visual} + w^{Audio}}.$$
 (22)

This fusion manner is termed **EEG-Visual-Audio-Fusion**. Here,  $y_i^{EEG}$ ,  $y_i^{Visual}$ , and  $y_i^{Audio}$  are the detected affective results by EEG, visual, and audio information, respectively.  $w^{EEG}$ ,  $w^{Visual}$ , and  $w^{Audio}$  are the corresponding fusion weight. For cross-comparison with EEG-Video-Fusion, the shared weights ( $w^{EEG}$ ,  $w^{Visual}$ ,  $w^{Audio}$ ) are set equally as well in EEG-Visual-Audio-Fusion.

We compare the affective detection performance of EEG-Video-Fusion and EEG-Visual-Audio-Fusion for valence and arousal on MAHNOB-HCI and DEAP databases, with different groundtruth types (aggregated and non-aggregated) and EEG models (MsDANN and MsNN). The corresponding results are reported in Table 10 and 11, which show EEG-Video-Fusion outperforms in most cases on the two databases. For example, compared to EEG-Visual-Audio-Fusion, EEG-Video-Fusion leads to 17.42% (Pacc) and 23.58% (Pf) increase for valence and 19.29% (Pacc) and 25.87% (Pf) increase for arousal on MAHNOB-HCI, when aggregated groundtruth is adopted with MsDANN. A similar increase trend is also observed for non-aggregated case on MAHNOB-HCI with MsNN. For DEAP database, compared to EEG-Visual-Audio-Fusion, it is observed that EEG-Video-Fusion performs better on valence (with increase rates of 17.26% for aggregated and 13.79% for non-aggregated in terms of  $P_f$ ), but performs similarly on arousal. One possible reason could be the elicited arousal largely relied on the audio content, when music videos were used as the emotion-evoking materials. The above comparison results demonstrate that the adopted EEG-Video-Fusion in the proposed EEG-AVE model achieves a better and more reliable cross-individual affective detection performance across emotion dimensions and databases.

## 4.5. Conclusion

In this paper, we propose a novel affective detection model (EEG-AVE) with an embedding protocol, where both EEG-based emotional individual preferences and audio-visual-based intrinsic emotions are incorporated to tackle the problem of the individual differences in EEG processing. The multimodal information is analyzed and compensated to realize efficient and effective EEG-based

#### Table 10

Affective detection performance with different fusion manners on MAHNOB-HCI and DEAP databases using MsDANN.

	Fusion Manner	MAHNOB	HCI			DEAP	DEAP			
		Aggregated		Non-Aggregated		Aggregated		Non-Aggregated		
		Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	
Valence	EEG-Video-Fusion	90.21	90.45	71.13	66.83	75.26	77.16	68.50	68.81	
	EEG-Visual-Audio-Fusion	76.83	73.19	70.54	58.49	70.45	65.80	68.75	60.47	
Arousal	EEG-Video-Fusion	85.59	86.55	66.47	63.25	71.92	79.50	54.52	60.80	
	EEG-Visual-Audio-Fusion	71.75	68.76	63.04	50.44	73.10	79.89	55.37	60.45	

#### Table 11

Affective detection performance with different fusion manners on MAHNOB-HCI and DEAP databases using MsNN.

	Fusion Manner	MAHNOB-	HCI			DEAP	DEAP					
		Aggregated		Non-Aggregated		Aggregated		Non-Aggregated				
		Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$	Pacc	$P_f$			
Valence	EEG-Video-Fusion	82.00	83.95	63.58	62.35	66.03	72.55	55.59	61.84			
	EEG-Visual-Audio-Fusion	72.78	73.88	59.05	58.94	70.46	70.52	61.45	58.48			
Arousal	EEG-Video-Fusion	78.39	81.42	58.00	58.17	72.62	80.07	50.62	57.59			
	EEG-Visual-Audio-Fusion	75.34	71.82	67.58	49.22	73.11	79.95	50.33	56.14			

affective detection. The experimental results show that the proposed EEG-AVE model achieves promising affective detection results, compared to the state-of-the-art methods. Besides, aiming at characterizing dynamic, informative, and domain-invariant EEG features across individuals, we develop a deep neural network with a transfer learning method (MsDANN) to solve the problem of the individual differences in the EEG data processing and investigate the performance variants with different neural network architectures (with or without domain adaptation). Our analysis demonstrates a superior cross-individual result is achieved under an evaluation of the leave-one-individual-out cross-validation individual-independent method. Furthermore, we utilize two well-known pretrained CNNs for semantic audio-visual feature extraction and introduce hypergraph theory to decode deep visual features, deep auditory features, and deep audio-visual fusion features for intrinsic emotions estimation. The possibility of affective detection using the multimedia materials is verified and the benefit of the proposed embedding strategy is examined. These results show both EEG signals and audio-visual information play important and helpful roles in affective detection, and the proposed EEG-AVE model could be applied to boost the development of affective brain-computer interface in real applications.

## **Declaration of Competing Interest**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

This work was supported in part by the National Natural Science Foundation of China under Grant 61906122, in part by Shenzhen– Hong Kong Institute of Brain Science-Shenzhen Fundamental Research Institutions (2021SHIBS0003), in part by Shenzhen Science and Technology Research and Development Fund for Sustainable Development Project (KCXFZ20201221173613036), in part by the Tencent "Rhinoceros Birds"-Scientific Research Foundation for Young Teachers of Shenzhen University, and in part by the High Level University Construction under Grant 000002110133.

## References

- T. Song, W. Zheng, P. Song, Z. Cui, Eeg emotion recognition using dynamical graph convolutional neural networks, IEEE Transactions on Affective Computing 11 (3) (2018) 532–541.
- [2] J. Li, S. Qiu, Y.-Y. Shen, C.-L. Liu, H. He, Multisource transfer learning for crosssubject eeg emotion recognition, IEEE transactions on cybernetics 50 (7) (2019) 3281–3293.
- [3] J. Li, S. Qiu, C. Du, Y. Wang, H. He, Domain adaptation for eeg emotion recognition based on latent representation similarity, IEEE Transactions on Cognitive and Developmental Systems 12 (2) (2019) 344–353.
- [4] Y. Cimtay, E. Ekmekcioglu, Investigating the use of pretrained convolutional neural network on cross-subject and cross-dataset eeg emotion recognition, Sensors 20 (7) (2020) 2034.
- [5] Y. Yin, X. Zheng, B. Hu, Y. Zhang, X. Cui, Eeg emotion recognition using fusion model of graph convolutional neural networks and lstm, Applied Soft Computing 100 (2021) 106954.
- [6] S. Jirayucharoensak, S. Pan-Ngum, P. Israsena, Eeg-based emotion recognition using deep learning network with principal component based covariate shift adaptation, The Scientific World Journal 2014 (2014).
- [7] W.-L. Zheng, B.-L. Lu, Investigating critical frequency bands and channels for eeg-based emotion recognition with deep neural networks, IEEE Transactions on Autonomous Mental Development 7 (3) (2015) 162–175.
- [8] H. Cui, A. Liu, X. Zhang, X. Chen, K. Wang, X. Chen, Eeg-based emotion recognition using an end-to-end regional-asymmetric convolutional neural network, Knowledge-Based Systems 205 (2020) 106243.
- [9] S. Liu, X. Wang, L. Zhao, B. Li, W. Hu, J. Yu, and Y. Zhang, "3dcann: A spatiotemporal convolution attention neural network for eeg emotion recognition," IEEE Journal of Biomedical and Health Informatics, 2021.
- [10] W.-L. Zheng, Y.-Q. Zhang, J.-Y. Zhu, B.-L. Lu, Transfer components between subjects for eeg-based emotion recognition, in: 2015 international conference on affective computing and intelligent interaction (ACII), IEEE, 2015, pp. 917– 922.
- [11] W.-L. Zheng, B.-L. Lu, Personalizing eeg-based affective models with transfer learning, in: Proceedings of the twenty-fifth international joint conference on artificial intelligence, 2016, pp. 2732–2738.
- [12] Y.-P. Lin, T.-P. Jung, Improving eeg-based emotion classification using conditional transfer learning, Frontiers in human neuroscience 11 (2017) 334.
- [13] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," The journal of machine learning research, vol. 17, no. 1, pp. 2096–2030, 2016.
- [14] Y. Li, W. Zheng, Y. Zong, Z. Cui, T. Zhang, X. Zhou, A bi-hemisphere domain adversarial neural network model for eeg emotion recognition, IEEE Transactions on Affective Computing (2018).
- [15] Y. Wang, L. Guan, A.N. Venetsanopoulos, Kernel cross-modal factor analysis for information fusion with application to bimodal emotion recognition, IEEE Transactions on Multimedia 14 (3) (2012) 597–607.
- [16] S. Mo, J. Niu, Y. Su, S.K. Das, A novel feature set for video emotion recognition, Neurocomputing 291 (2018) 11–20.
- [17] S. Zhang, S. Zhang, T. Huang, W. Gao, Q. Tian, Learning affective features with a hybrid deep model for audio-visual emotion recognition, IEEE Transactions on Circuits and Systems for Video Technology 28 (10) (2017) 3030–3043.

## Z. Liang, X. Zhang, R. Zhou et al.

- [18] B. Liang, H. Su, R. Yin, L. Gui, M. Yang, Q. Zhao, X. Yu, R. Xu, Beta distribution guided aspect-aware graph for aspect category sentiment analysis with affective knowledge, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, 2021, pp. 208–218.
- [19] E. Acar, F. Hopfgartner, and S. Albayrak, "A comprehensive study on mid-level representation and ensemble learning for emotional analysis of video material," Multimedia Tools and Applications, vol. 76, no. 9, pp. 11 809–11 837, 2017.
- [20] J. Cheng, M. Chen, C. Li, Y. Liu, R. Song, A. Liu, X. Chen, Emotion recognition from multi-channel eeg via deep forest, IEEE Journal of Biomedical and Health Informatics 25 (2) (2020) 453–464.
- [21] S. Kim, H.-J. Yang, N.A.T. Nguyen, S.K. Prabhakar, and S.-W. Lee, "Wedea: A new eeg-based framework for emotion recognition," IEEE Journal of Biomedical and Health Informatics, 2021.
- [22] F. Noroozi, M. Marjanovic, A. Njegus, S. Escalera, G. Anbarjafari, Audio-visual emotion recognition in video clips, IEEE Transactions on Affective Computing 10 (1) (2017) 60–75.
- [23] S. Poria, E. Cambria, R. Bajpai, A. Hussain, A review of affective computing: From unimodal analysis to multimodal fusion, Information Fusion 37 (2017) 98–125.
- [24] S. Wang, S. Chen, Q. Ji, Content-based video emotion tagging augmented by users' multiple physiological responses, IEEE Transactions on Affective Computing 10 (2) (2017) 155–166.
- [25] Y. Tonoyan, D. Looney, D.P. Mandic, M.M. Van Hulle, Discriminating multiple emotional states from eeg using a data-adaptive, multiscale informationtheoretic approach, International journal of neural systems 26 (02) (2016) 1650005.
- [26] K. Michalopoulos and N. Bourbakis, "Application of multiscale entropy on eeg signals for emotion detection," in 2017 IEEE EMBS International Conference on Biomedical & Health Informatics (BHI). IEEE, 2017, pp. 341–344.
- [27] A. Martínez-Rodrigo, B. García-Martínez, R. Alcaraz, P. González, A. Fernández-Caballero, Multiscale entropy analysis for recognition of visually elicited negative stress from eeg recordings, International journal of neural systems 29 (02) (2019) 1850038.
- [28] K. Simonyan and A. Zisserman, "Very deep convolutional networks for largescale image recognition," arXiv preprint arXiv:1409.1556, 2014.
- [29] S. Hershey, S. Chaudhuri, D.P. Ellis, J.F. Gemmeke, A. Jansen, R.C. Moore, M. Plakal, D. Platt, R.A. Saurous, B. Seybold, et al., Cnn architectures for large-scale audio classification, in: 2017 ieee international conference on acoustics, speech and signal processing (icassp), IEEE, 2017, pp. 131–135.
- [30] A. Sengupta, Y. Ye, R. Wang, C. Liu, K. Roy, Going deeper in spiking neural networks: Vgg and residual architectures, Frontiers in neuroscience 13 (2019) 95.
- [31] W. Han, T. Jiang, Y. Li, B. Schuller, and H. Ruan, "Ordinal learning for emotion recognition in customer service calls," in ICASSP 2020–2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2020, pp. 6494–6498.
- [32] R.-N. Duan, J.-Y. Zhu, and B.-L. Lu, "Differential entropy feature for eeg-based emotion classification," in 2013 6th International IEEE/EMBS Conference on Neural Engineering (NER). IEEE, 2013, pp. 81–84.
- [33] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, Advances in neural information processing systems 25 (2012) 1097–1105.
- [34] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 1–9.
- [35] L. Shi, K. Du, C. Zhang, H. Ma, and W. Yan, "Lung sound recognition algorithm based on vggish-bigru," IEEE Access, vol. 7, pp. 139 438–139 449, 2019.
- [36] S. Kurada, A. Kurada, Poster: Vggish embeddings based audio classifiers to improve parkinson's disease diagnosis, in: 2020 IEEE/ACM International Conference on Connected Health: Applications, Systems and Engineering Technologies (CHASE), ACM, 2020, pp. 9–11.
- [37] A. Ducournau, S. Rital, A. Bretto, B. Laget, "A multilevel spectral hypergraph partitioning approach for color image segmentation," in: 2009 IEEE International Conference on Signal and Image Processing Applications. IEEE, 2009, pp. 419–424.
- [38] D. Zhou, J. Huang, B. Schölkopf, Learning with hypergraphs: Clustering, classification, and embedding, Advances in neural information processing systems 19 (2006) 1601–1608.
- [39] M. Soleymani, J. Lichtenauer, T. Pun, M. Pantic, A multimodal database for affect recognition and implicit tagging, IEEE transactions on affective computing 3 (1) (2011) 42–55.
- [40] S. Koelstra, C. Muhl, M. Soleymani, J.-S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, I. Patras, Deap: A database for emotion analysis; using physiological signals, IEEE transactions on affective computing 3 (1) (2011) 18–31.
- [41] Y. Zhu, S. Wang, Q. Ji, Emotion recognition from users' eeg signals with the help of stimulus videos, in: 2014 IEEE international conference on multimedia and expo (ICME), IEEE, 2014, pp. 1–6.
- [42] X. Huang, J. Kortelainen, G. Zhao, X. Li, A. Moilanen, T. Seppänen, M. Pietikäinen, Multi-modal emotion analysis from facial expressions and electroencephalogram, Computer Vision and Image Understanding 147 (2016) 114–124.

- [43] S. Rayatdoost, M. Soleymani, Cross-corpus eeg-based emotion recognition, in: 2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP), IEEE, 2018, pp. 1–6.
- [44] Z. Liang, S. Oba, S. Ishii, An unsupervised eeg decoding system for human emotion recognition, Neural Networks 116 (2019) 257–268.
- [45] M. Soleymani, G. Chanel, J.J. Kierkels, T. Pun, Affective ranking of movie scenes using physiological signals and content analysis, in: Proceedings of the 2nd ACM Workshop on Multimedia Semantics, 2008, pp. 32–39.
- [46] M. Soleymani, J.J. Kierkels, G. Chanel, and T. Pun, "A bayesian framework for video affective representation," in: 2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops. IEEE, 2009, pp. 1–7.
- [47] A. Yazdani, K. Kappeler, T. Ebrahimi, Affective content analysis of music video clips, in: Proceedings of the 1st international ACM workshop on Music information retrieval with user-centered and multimodal strategies, 2011, pp. 7–12.



**Zhen Liang** received her Ph.D. degree from The Hong Kong Polytechnic University, Hong Kong, in 2013. From 2012 to 2017, she was an algorithm development scientist at NeuroSky, Inc., Hong Kong. From 2018 to 2019, she was a specially-appointed assistant professor at Graduate School of Informatics, Kyoto University, Japan. She is currently an assistant professor in the School of Biomedical Engineering, Health Science Center, Shenzhen University, China. Her current research interests include brain encoding and decoding systems, affective computing, visual attention, and neural engineering.



Xihao Zhang is a master student in the School of Biomedical Engineering, Health Science Center, Shenzhen University, China. His research interests include affective computing and multimedia processing.



**Rushuang Zhou** is a bachelor student in the School of Biomedical Engineering, Health Science Center, Shenzhen University, China. His research interests include affective computing and transfer learning.



Li Zhang received his Ph.D. degree from the Department of Electrical and Electronics, University of Hong Kong in 2017 and joined School of Biomedical Engineering, Health Science Center, Shenzhen University in 2018 as an associate researcher. His current research interests mainly focus on biomedical signal processing, numerical optimization, and imaging genetics.

## Neurocomputing 510 (2022) 107-121





**Linling Li** received her Ph. D. degree from the University of Chinese Academy of Sciences, China (2017) and then received her postdoctoral training at University of Shenzhen, China. She is now an associate researcher in School of Biomedical Engineering, Health Science Center, Shenzhen University, China. Her current research interests focus on the development of EEG neurofeedback technique and related neuroimaging mechanisms.



**Zhiguo Zhang** received his B.Eng. degree from Tianjin University in 2000, his M.Phil. degree from the University of Science and Technology of China in 2003, and his Ph.D. degree from the University of Hong Kong in 2008. He is now a professor with the School of Biomedical Engineering, Health Science Center, Shenzhen University, China. His research interests include biomedical signal processing, neural engineering, and braininspired computation.



**Gan Huang** received his Ph.D. degree from Shanghai Jiao Tong University, China (2013) and then received his postdoctoral training at Universit? Catholique de Louvain in Belgium. He is now an assistant professor in School of Biomedical Engineering, Shenzhen University, China. He has published several high-quality papers in journals including Neuroimage, Neural Network, and Neurocomputing. His current research interests focus on neural modulation, brain-computer interface, and neuroprosthesis.