SPATIAL-TEMPORAL TRANSFORMERS FOR EEG EMOTION RECOGNITION

Jiyao Liu, Hao Wu, Li Zhang, Yanxi Zhao

School of Computer Science, Northwestern Polytechnical University, Xi'an, China

ABSTRACT

Electroencephalography (EEG) is a popular and effective tool for emotion recognition. However, the propagation mechanisms of EEG in the human brain and its intrinsic correlation with emotions are still obscure to researchers. This work proposes four variant transformer frameworks (spatial attention, temporal attention, sequential spatial-temporal attention and simultaneous spatial-temporal attention) for EEG emotion recognition to explore the relationship between emotion and spatial-temporal EEG features. Specifically, spatial attention and temporal attention are to learn the topological structure information and time-varying EEG characteristics for emotion recognition respectively. Sequential spatial-temporal attention does the spatial attention within a one-second segment and temporal attention within one sample sequentially to explore the influence degree of emotional stimulation on EEG signals of diverse EEG electrodes in the same temporal segment. The simultaneous spatial-temporal attention, whose spatial and temporal attention are performed simultaneously, is used to model the relationship between different spatial features in different time segments. The experimental results demonstrate that simultaneous spatial-temporal attention leads to the best emotion recognition accuracy among the design choices, indicating modeling the correlation of spatial and temporal features of EEG signals is significant to emotion recognition.

Index Terms— EEG, emotion recognition, transformer

1. INTRODUCTION

EEG emotion recognition is to detect the current emotional states of the subjects [1], [2], [3]. In recent years, with the development of deep learning and the availability of EEG data, many emotion recognition methods based on neural networks have dominated the state-of-art position [4, 5]. In general, the EEG signals collected by a spherical EEG cap have three-dimensional characteristics which are spatial, spectral and temporal respectively. Many researchers have drawn attention to how to effectively utilize time-varying spatial and temporal features from multi-channel brain signals.

In order to model the space relationships among multichannel EEG signals, a hierarchical convolutional neural network (CNN) is proposed by Li et al. [5] to capture spatial information among different channels. A deep CNN model is present by Zhang et al. [6] to capture the spatio-temporal robust feature representation of the raw EEG data stream for motion intention classification. A utilization of multi-layer CNN with no full connection layers is proposed by Lawhern et al. [7] for P300-based oddball recognition task, finger motor task and motor imagination task.

Considering the change of EEG signals over time, an Echo State Network (ESN) is present by Fourati et al. [8], ESN used recursive layer to map the EEG signal into the high-dimension state space. A two-layer long short term memory (LSTM), which uses EEG signal as the input, is adopted by Alhagry et al. [9] and obtain promising EEG emotion classification results. A deep recursive convolutional neural network (R-CNN) is present by Bashivan et al. [10], the proposed R-CNN gets a satisfactory result on the task mental load classification based on EEG signal.

Most of the above works are on the basis of convolution or recursive operation. CNN is good at modeling local receptive field message, while pays less attention to the global information. Recurrent Neural Networks (RNN) network is relatively weak to capture the spatial information and its parallel computational efficiency is slower. To solve the above weaknesses, some works lead attention mechanism into CNN and RNN.

Since different spatial-temporal features have different contributions to emotion recognition, they should be assigned to different weight in the classifier recognizing emotions. A LSTM with attention mechanism is proposed by Kim et al. [11], the network assigns weights to the emotional states appearing at specific moments to conduct two-level and three-level classification on the valence and arousal emotion models. A fresh multi-channel model on the basis of sparse graphic attention long short term memory (SGA-LSTM) is present by Liu et al. [12] to classify EEG emotion.

As is mentioned above, existing works have attained gratifying results. However, The transmission characteristic as well as spatial-temporal relevance of different EEG electrodes are more or less neglected in most of them. The changeless size kernels in convolution operation [13] may damage the spatial correlation of EEG signals. Though the RNN operation [14] takes the temporal features of EEG signals into

consideration, it ignores the spatial relation among EEG electrodes. Furthermore, on account of the diverse impedance of various brain areas, there may be a slight error in time between the EEG signal displayed by the EEG collection device and the real EEG signal, that is, EEG signals may delay varies with different EEG electrodes.

To deal with the mentioned issues, we present a fresh EEG emotion transformer (EeT) framework built exclusively on self-attention blocks. The variants of self-attention block include spatial (S) attention, temporal (T) attention, sequential spatial-temporal (S-T) attention and simultaneous spatialtemporal (S+T) attention. The spatial attention is to learn the spatial structure information. The temporal attention is to learn the correlation between EEG signals and emotional stimuli as well as temporal changes. The sequential spatialtemporal attention is to do spatial attention within the same time segment and temporal attention among different time segments in one sample. The simultaneous spatial-temporal attention is to do the two attention simultaneously. Experimental setups are elaborately picked to study the effects of spatial and temporal EEG signals on emotion recognition, and whether there is some correlation between the features of different channels at different time segments.

2. VARIANTS OF EEG EMOTION TRANSFORMERS

2.1. Framework of EeT

EEG-based emotion recognition is to classify the emotion states according to the EEG signal. As illustrated in Fig. 1, the overview of EeT includes four modules, namely the feature preparation module, the spatial-temporal attention module, deep neural network (DNN) module and classification module. We focus on the design of self-attention module which includes spatial attention, temporal attention, sequential spatial-temporal attention and simultaneous spatial-temporal attention.

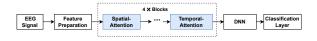


Fig. 1: Overview of EEG Emotion Recognition Transformer (EeT)

Emotion recognition based on EEG is to learn a function *f* which maps the raw signals to emotion tags:

$$Y = f(X') \tag{1}$$

where X' represents the EEG features. f represents the mapping function i.e., convolutional neural network transformation. $Y \in \{y_1, y_2, ..., y_n\}$ represents the emotional tags. In our work, the cross entropy is adopted as the loss function, which can be defined as:

$$L = -\sum_{c=1}^{C} y_c log(y'_c)$$
 (2)

where L denotes the loss function of the EEG emotion classification, C denotes the number of emotion states, y_c is the ground-truth emotion tag and $y_c^{'}$ is the predictor of neural networks.

2.2. Preprocessing

The EEG features can be denoted as $X = (F_1, F_2, ..., F_T) \in$ $\mathbb{R}^{C \times S \times T}$, where C is 5, equals to the number of frequency bands (δ [1-4Hz], θ [4-8Hz], α [8-14Hz], β [14-31Hz], and γ [31-50 Hz]) used to compute the EEG features. S equals to the number of electrodes in the EEG cap and T equals to the number of the time slots in a EEG sample. $F_t = (B_1, B_2, ..., B_S) \in \mathbb{R}^{C \times S} (t \in \{1, 2, ..., T\})$ is the one-second EEG feature. $B_s = (b_1, b_2, ..., b_C) \in$ $\mathbb{R}^C(s \in \{1, 2, ..., S\})$ presents the feature of one EEG channel. Specifically, we map the S electrodes into a $V \times H$ matrix according to the layout of the EEG cap in order to preserve the spatial topology structure information of the EEG cap, then we take advantage of the linear interpolation [15] to replenish the spatial information which are not collected by EEG acquisition equipment. The re-organized feature can be represent as $F_t^{'}=(B_1,B_2,...,B_{V\times H})\in\mathbb{R}^{C\times V\times H}$ $(t\in 1,2,...,T)$ for a one-second EEG slice and $X' = (F_1, F_2, ..., F_T) \in \mathbb{R}^{C \times V \times H \times T}$ for the whole EEG sample.

2.3. Positional Encoding for Spatial EEG Electrodes

We divided the EEG feature of each second $F_t(i=1,2,3..S)$ into G non-overlapping regions, just like the different brain regions in neuroscience. Here we regroup the $V \times H$ matrices into region sequences, the size of each divided region is $P \times P$, so we get $G = VH/P^2$ regions. Each region is flatten into a vector $I(x)_{(p,t)} \in \mathbb{R}^{5P^2}$ with p=1,2...,G representing spatial layout of EEG electrodes and t=1,2,...T denoting the index over seconds. Then we linearly map each region $I(x)_{(p,t)}$ into a latent vector $z_{(p,t)}^{(0)} \in \mathbb{R}^D$ by means of learnable matrix $M \in \mathbb{R}^{D \times 5P^2}$:

$$z_{(p,t)}^{(0)} = M \otimes I(x)_{(p,t)} + e_{(p,t)}^{position}$$
 (3)

where \otimes is matrix multiplication and $e_{(p,t)}^{pos} \in \mathbb{R}^D$ stands for a learnable position embedding to encode the spatial-temporal position of each brain region. The resulting sequence of embedding vectors $z_{(p,t)}^{(0)}$ stands for the input to the next layer of the self-attention block. Note that z^i is output of the ith layer in self-attention block. p=1,...,G and t=1,...,T are the spatial locations and indexes over time segments respectively.

2.4. Query-Key-Value Mechanism

Our Transformer consists of L encoding blocks. Instead of performing a single attention function, we use different pro-

jected versions of queries, keys and values to perform the attention function in parallel, which is called multi-head attention. At each block l, the query/key/value vectors are computed for each region from the representation $z_{(n,t)}^{(l-1)}$ encoded by the preceding block:

$$q_{(p,t)}^{(l,a)} = W_Q^{(l,a)} z_{(p,t)}^{(l-1)} \in \mathbb{R}^{D_h}$$
(4)

$$k_{(p,t)}^{(l,a)} = W_K^{(l,a)} z_{(p,t)}^{(l-1)} \in \mathbb{R}^{D_h}$$

$$v_{(p,t)}^{(l,a)} = W_V^{(l,a)} z_{(p,t)}^{(l-1)} \in \mathbb{R}^{D_h}$$
(6)

$$v_{(p,t)}^{(l,a)} = W_V^{(l,a)} z_{(p,t)}^{(l-1)} \in \mathbb{R}^{D_h}$$
 (6)

 $z_{(p,t)}^{(l-1)},$ which is the output of the previous block, need to be layer normalized before the above operations. a=1,2...,Adenotes an index over multiple attention heads and A is the number of attention heads.

2.5. Variants of Attention Mask Learning

The variants of self-attention block include spatial attention (S), temporal attention (T), sequential spatial-temporal attention (S-T) and simultaneous spatial-temporal attention (S+T). The spatial attention is to learn the spatial structure information while the temporal attention is to the relationship between EEG and time. The sequential spatial-temporal attention is the concatenation of two operations. The simultaneous spatial-temporal attention is to do the two operations simultaneously.

2.5.1. Spatial Attention

In the case of spatial attention, the self-attention weights $\alpha_{(p,t)}^{(a,l)} \in \mathbb{R}^{N+1}$ for query brain region (p,t) are given by:

$$\alpha_{(p,t)}^{(l,a) \text{ spatial}} = \sigma \left(\frac{q_{(p,t)}^{(l,a)^{\top}}}{\sqrt{D_h}} \cdot \left[k_{(0,0)}^{(l,a)} \left\{ k_{(p',t)}^{(l,a)} \right\}_{p'} \right] \right)$$
(7)

where p' denotes the index of the brain regions. σ denotes the softmax activation function. The formula is to consider that different brain regions react differently under the same emotional stimulation thus different weights are given to the features of different brain regions.

2.5.2. Temporal Attention

For the temporal attention, the self-attention weights $\alpha_{(p,t)}^{(l,a)} \in$ \mathbb{R}^{T+1} for query brain region (p,t) are given by:

$$\alpha_{(p,t)}^{(l,a) \text{ temporal}} = \sigma \left(\frac{q_{(p,t)}^{(l,a)^{\top}}}{\sqrt{D_h}} \cdot \left[k_{(0,0)}^{(l,a)} \left\{ k_{(p,t')}^{(l,a)} \right\}_{t'} \right] \right)$$
(8)

where t' denotes the index of the time slots. In this formula, different weights are given to the of different time slots in consideration of the change of EEG signals with emotional stimulation and time.

2.5.3. Sequential Spatial-Temporal (S-T) Attention

The sequential spatial-temporal attention is to do spatial attention within the same time segment and temporal attention among different time segments. Firstly, the spatial selfattention weights are calculated as Eq. 7. Then the temporal attention weights are learned by Eq. 8 from the output of spatial attention layer. (S-T) Attention comprehensively considers the attention of space and time, but the default is that the spatial features in the same time period are closely related, while the features of different brain in different time are weakly related.

2.5.4. Simultaneous Spatial-Temporal (S+T) Attention

The simultaneous spatial-temporal (S+T) attention is to do spatial and temporal attention simultaneously. The selfattention weights $\alpha_{(p,t)}^{(a,l)} \in \mathbb{R}^{NT+1}$ for query brain region (p,t) are given by:

$$\alpha_{(p,t)}^{(l,a)} = \sigma \left(\frac{q_{(p,t)}^{(l,a)}}{\sqrt{D_h}} \right) \cdot \left[k_{(0,0)}^{(l,a)} \left\{ k_{(p',t')}^{(l,a)} \right\} p' = 1, ..., N \right] t' = 1, ..., T \right]$$

$$(9)$$

Different from S-T Attention, which regards space and time separately, the S+T attention considers that the spatial information in the same time point and different time points are both strongly correlated.

2.6. Multi-head Attention Recalibration

The encoding $z_{(p,t)}^{(l)}$ at block l is obtained by the first computing the weighted sum of value vectors using self-attention coefficients from each attention head:

$$s_{(p,t)}^{(l,a)} = \alpha_{(p,t),(0,0)}^{(l,a)} v_{(0,0)}^{(l,a)} + \sum_{p'=1}^{N} \sum_{t'=1}^{F} \alpha_{(p,t),(p',t')}^{(l,a)} v_{(p',t')}^{(l,a)},$$
(10)

As is mentioned above, a denotes an index over multiple attention heads and l is the index of the blocks. Then, the concatenation of these vectors from all heads is projected and passed through an MLP, using residual connections after each operation:

$$z'_{(p,t)}^{l} = W_{O}^{(l-1)} \begin{bmatrix} s_{(p,t)}^{(l,1)} \\ \vdots \\ \vdots \\ s_{(p,t)}^{(l,A)} \end{bmatrix} + z_{(p,t)}^{(l-1)}, \tag{11}$$

where W_O is the **Value** of $z_{(p,t)}^{(l-1)}$ by concatenating $v_{(p,t)}^{(l,a)}$.

$$z_{(p,t)}^{l} = MLP(z_{(p,t)}^{\prime l}) + z_{(p,t)}^{\prime l}$$
(12)

The $z'^{l}_{(p,t)}$ goes through the MLP layer to get the output of

3. EXPERIMENTS

3.1. Datasets

We validate our model on SEED [16, 17], SEED-IV [18] and Deap [19] databases.

Deap dataset is a open source dataset including diverse physiological signals with emotion evaluations provided by the research team of Queen Mary University in London. It records the EEG, ECG, EMG and other bioelectrical signals of 32 subjects induced by watching 40 one-minute music videos of different emotional tendencies. The subjects evaluated the videos' emotion categories on scale of one to nine in dimension of arousal, valence, liking, dominance and familiarity. Valence reports the degree of subjects' joy, the greater the valence value, the higher the joy degree. Arousal reports the subjects' emotional intensity, the higher the arousal value, the more intense and perceptible the emotion. The rating value from small to large indicates the emotion metric is from negative to positive or from weak to strong. The 40 stimulus videos include 20 high valence/arousal stimuli and 20 low valence/arousal stimuli.

SEED contains three different categories of emotion, namely positive, negative, and neutral. Fifteen participants' EEG data of the dataset were collected while they were watching the stimulus videos. The videos are carefully selected and can elicit a single desired target emotion. With an interval of about one week, each subject participated in three experiments, and each session contained 15 film clips. The participants are asked to give feedback immediately after each experiment. The EEG signals of 62 channels are recorded at a sampling frequency of 1000 Hz and down-sampled with 200 Hz. The Differential entropy [17] DE features are precomputed over different frequency bands for each sample in each channel.

SEED-IV contains four different categories of emotions, including happy, sad, fear, and neutral emotion. The experiment consists of 15 participants. Three experiments are designed for each participant on different days, and each session contains 24 video clips and six clips for each type of emotion in each session. After each experiment, the subjects are asked to give feedback, while 62 EEG signals of the subjects are recorded. The EEG signals are sliced into 4-second non-overlapping segments and down-sampled with 128 Hz. The DE feature is also pre-computed over five frequency bands in each channel.

3.2. Experimental Setup

We train our model on NVIDIA RTX 2080 GPU. Cross entropy loss is used as the loss function. The optimizer is Adam. The initial learning rate is set to 1e-3 with multi-step decay to 1e-7. The number of the attention blocks is set to 4 and the length of each sample is set to 10s. We conduct experiments on each subject. For each experiment, we randomly shuffle

the samples and use 5-fold cross validation. The ratio of the training set to test set is 9:6.

3.3. Compared Models

We compare the proposed EeT with the following competitive models.

SVM [20] is a least squares support vector machine classifier. DBN [21] is deep Belief Networks investigate the critical frequency bands and channels. DGCNN [22] is Dynamical Graph Convolutional Neural Networks model the multichannel EEG features. BiDANN [23] is bi-hemispheres domain adversarial neural network maps the EEG feature of both hemispheres into discriminative feature spaces separately. BiHDM [24] is bi-hemispheric discrepancy model learns the asymmetric differences between two hemispheres for EEG emotion recognition. 3D-CNN with PST-Attention [25] is a self-attention module combined with 3D-CNN to learn critical information among different dimensions of EEG feature. LSTM [9] is a time series model for identifying continuous dimension emotion. 3D-CNN [26] is 3D-CNN model to recognize arousal and valence. BT [27] is deep convolution neural network for continuous dimension emotion recogni-

3.4. Experimental Results and Analysis

Table 1: Experimental Results on DEAP Dataset

| Models | Arousal | | Valence | |
|----------------------------|---------|--------|---------|--------|
| | acc(%) | F1 | acc(%) | F1 |
| LSTM[9] | 85.65 | - | 85.45 | - |
| 3D-CNN [26] | 88.49 | - | 87.44 | - |
| BT[27] | 86.18 | - | 86.31 | - |
| $EeT \sim (S+T Attention)$ | 93.34 | 0.9326 | 92.86 | 0.9196 |

Table 1 presents the average accuracy (acc) and F1 value (F1) of the compared models for EEG based emotion recognition on the DEAP datasets. Compared with 3D-CNN [26], the acc of the proposed simultaneous spatial-temporal (S+T) attention EeT framework have 4.85%/5.42% improvement. EeT with S+T attention also gets superior performance compared with other competitive models.

Table 2: Experimental Results on SEED and SEED-IV Dataset

| Models | SEED | | SEED-IV | |
|-------------------------------|----------|---------|----------|---------|
| Wiodels | Mean (%) | Std (%) | Mean (%) | Std (%) |
| SVM [20] | 83.99 | 9.72 | 56.61 | 20.05 |
| DBN [21] | 86.08 | 8.34 | 66.77 | 7.38 |
| DGCNN [22] | 90.40 | 8.49 | 69.88 | 16.29 |
| BiDANN [23] | 92.38 | 7.04 | 70.29 | 12.63 |
| BiHDM [24] | 93.12 | 6.06 | 74.35 | 14.09 |
| 3D-CNN with PST-Attention[25] | 95.76 | 4.98 | 82.73 | 8.96 |
| EeT (S+T Attention) | 96.28 | 4.39 | 83.27 | 8.37 |

Table 2 presents the average accuracy (Mean) and standard deviation (Std) of the compared models for EEG based

emotion recognition on SEED and SEED-IV datasets. Compared with 3D-CNN with PST-Attention, the means of the proposed joint spatial+temporal (S+T) attention EeT framework have 0.52%/0.54% improvements on SEED and SEED-IV. The Stds of Eet with S+T attention achieve 0.59%/0.59% reductions on SEED and SEED-IV respectively compared with those of 3D-CNN with PST-Attention. Moreover, EeT with S+T attention gets superior performance compared with other competitive models.

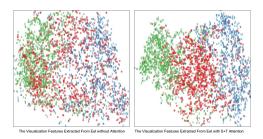


Fig. 2: The Visualization of High Level Features Extracted From Eet

We use t-SNE [28] to visualize the high-level bottleneck features from the well trained Eet. As shown in Fig. 2, different colors represent different emotional labels, the distances of different classes in the high-level feature space of Eet with S+T Attention (the right part) are more dispersed than that of EeT without attention (the left part), which demonstrates that the high-level features learned with simultaneous spatial-temporal attention is more discriminative.

3.5. Ablation Experiments

Table 3: Experimental Results of Variant Transformers

| Models | SEED | | SEED-IV | | |
|---------------|----------|---------|----------|---------|--|
| | Mean (%) | Std (%) | Mean (%) | Std (%) | |
| S Attention | 93.14 | 9.31 | 73.31 | 13.67 | |
| T Attention | 92.74 | 10.21 | 72.37 | 12.83 | |
| S-T Attention | 95.65 | 6.73 | 80.31 | 8.51 | |
| S+T Attention | 96.28 | 4.39 | 83.27 | 8.37 | |

Table 3 presents the results of different variants of the proposed transformer framework, from which we can see that Joint Spatial-Temporal Attention gets the best results, achieving 0.63%/2.96% improvements compared with the second best variant, (S-T) Attention on SEED and SEED-IV respectively, indicating comprehensively considering the temporal and spatial characteristics of EEG may boost the emotion recognition results most notably. As for single dimensional attention, the results are a bit of lower than those of combined variants'. Spatial Attention is 0.4% /0.94% higher than that of Temporal Attention, implying the spatial dimension may have more emotion-related message than the temporal dimension.

4. CONCLUSION

In this paper, we propose a new EEG emotion recognition framework based on self-attention, which is built exclusively on self-attention. Our approach considers the relationship between emotion and brain regions, time series change as well as the intrinsic spatiotemporal characteristics of EEG signals. The results of our methods show that the attention mechanism can boost the performance of emotion recognition evidently. Furthermore, the simultaneous spatio-temporal attention gets the best results among the four designed structures, the result is also better than most state of the art methods, indicating that considering the spatio-temporal feature jointly and simultaneously is more in line with the transmission law of EEG signals in the human brain.

5. REFERENCES

- [1] Danny Oude Bos et al., "Eeg-based emotion recognition," *The influence of visual and auditory stimuli*, vol. 56, no. 3, pp. 1–17, 2006.
- [2] Roddy Cowie, Ellen Douglas-Cowie, Nicolas Tsapatsoulis, George Votsis, Stefanos Kollias, Winfried Fellenz, and John G Taylor, "Emotion recognition in human-computer interaction," *IEEE Signal processing* magazine, vol. 18, no. 1, pp. 32–80, 2001.
- [3] Hatice Gunes, Björn Schuller, Maja Pantic, and Roddy Cowie, "Emotion representation, analysis and synthesis in continuous space: A survey," in 2011 IEEE International Conference on Automatic Face & Gesture Recognition (FG). IEEE, 2011, pp. 827–834.
- [4] Abeer Al-Nafjan, Manar Hosny, Areej Al-Wabil, and Yousef Al-Ohali, "Classification of human emotions from electroencephalogram (eeg) signal using deep neural network," *Int. J. Adv. Comput. Sci. Appl*, vol. 8, no. 9, pp. 419–425, 2017.
- [5] Jinpeng Li, Zhaoxiang Zhang, and Huiguang He, "Hierarchical convolutional neural networks for eeg-based emotion recognition," *Cognitive Computation*, vol. 10, no. 2, pp. 368–380, 2018.
- [6] Dalin Zhang, Lina Yao, Xiang Zhang, Sen Wang, Weitong Chen, and Robert Boots, "Eeg-based intention recognition from spatio-temporal representations via cascade and parallel convolutional recurrent neural networks," arXiv preprint arXiv:1708.06578, 2017.
- [7] Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance, "Eegnet: a compact convolutional neural network for eeg-based brain-computer interfaces," *Journal of neural engineering*, vol. 15, no. 5, pp. 056013, 2018.

- [8] Rahma Fourati, Boudour Ammar, Chaouki Aouiti, Javier Sanchez-Medina, and Adel M Alimi, "Optimized echo state network with intrinsic plasticity for eeg-based emotion recognition," in *International Con*ference on Neural Information Processing. Springer, 2017, pp. 718–727.
- [9] Salma Alhagry, Aly Aly Fahmy, and Reda A El-Khoribi, "Emotion recognition based on eeg using 1stm recurrent neural network," *Emotion*, vol. 8, no. 10, pp. 355–358, 2017.
- [10] Pouya Bashivan, Irina Rish, Mohammed Yeasin, and Noel Codella, "Learning representations from eeg with deep recurrent-convolutional neural networks," *arXiv* preprint arXiv:1511.06448, 2015.
- [11] Youmin Kim and Ahyoung Choi, "Eeg-based emotion classification using long short-term memory network with attention mechanism," *Sensors*, vol. 20, no. 23, pp. 6727, 2020.
- [12] Suyuan Liu, Wenming Zheng, Tengfei Song, and Yuan Zong, "Sparse graphic attention lstm for eeg emotion recognition," in *International Conference on Neural Information Processing*. Springer, 2019, pp. 690–697.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
- [14] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [15] Heidi A Schlitt, L Heller, R Aaron, E Best, and DM Ranken, "Evaluation of boundary element methods for the eeg forward problem: effect of linear interpolation," *IEEE transactions on biomedical engineering*, vol. 42, no. 1, pp. 52–58, 1995.
- [16] Wei-Long Zheng and Bao-Liang Lu, "Investigating critical frequency bands and channels for eeg-based emotion recognition with deep neural networks," *IEEE Transactions on Autonomous Mental Development*, vol. 7, no. 3, pp. 162–175, 2015.
- [17] Ruo-Nan Duan, Jia-Yi Zhu, and Bao-Liang 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.
- [18] Wei-Long Zheng, Wei Liu, Yifei Lu, Bao-Liang Lu, and Andrzej Cichocki, "Emotionmeter: A multimodal framework for recognizing human emotions," *IEEE transactions on cybernetics*, vol. 49, no. 3, pp. 1110–1122, 2018.

- [19] Sander Koelstra, Christian Muhl, Mohammad Soleymani, Jong-Seok Lee, Ashkan Yazdani, Touradj Ebrahimi, Thierry Pun, Anton Nijholt, and Ioannis Patras, "Deap: A database for emotion analysis; using physiological signals," *IEEE transactions on affective computing*, vol. 3, no. 1, pp. 18–31, 2011.
- [20] Johan AK Suykens and Joos Vandewalle, "Least squares support vector machine classifiers," *Neural processing letters*, vol. 9, no. 3, pp. 293–300, 1999.
- [21] Wei-Long Zheng, Jia-Yi Zhu, Yong Peng, and Bao-Liang Lu, "Eeg-based emotion classification using deep belief networks," in 2014 IEEE International Conference on Multimedia and Expo (ICME). IEEE, 2014, pp. 1–6.
- [22] Tengfei Song, Wenming Zheng, Peng Song, and Zhen Cui, "Eeg emotion recognition using dynamical graph convolutional neural networks," *IEEE Transactions on Affective Computing*, vol. 11, no. 3, pp. 532–541, 2018.
- [23] Yang Li, Wenming Zheng, Zhen Cui, Tong Zhang, and Yuan Zong, "A novel neural network model based on cerebral hemispheric asymmetry for eeg emotion recognition.," in *IJCAI*, 2018, pp. 1561–1567.
- [24] Yang Li, Lei Wang, Wenming Zheng, Yuan Zong, Lei Qi, Zhen Cui, Tong Zhang, and Tengfei Song, "A novel bi-hemispheric discrepancy model for eeg emotion recognition," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 13, no. 2, pp. 354–367, 2020.
- [25] Jiyao Liu, Yanxi Zhao, Hao Wu, and Dongmei Jiang, "Positional-spectral-temporal attention in 3d convolutional neural networks for eeg emotion recognition," *Proc. APSIPA*, 2021.
- [26] Elham S Salama, Reda A El-Khoribi, Mahmoud E Shoman, and Mohamed A Wahby Shalaby, "Eeg-based emotion recognition using 3d convolutional neural networks," *Int. J. Adv. Comput. Sci. Appl*, vol. 9, no. 8, pp. 329–337, 2018.
- [27] JX Chen, PW Zhang, ZJ Mao, YF Huang, DM Jiang, and YN Zhang, "Accurate eeg-based emotion recognition on combined features using deep convolutional neural networks," *IEEE Access*, vol. 7, pp. 44317–44328, 2019.
- [28] Laurens Van der Maaten and Geoffrey Hinton, "Visualizing data using t-sne.," *Journal of machine learning research*, vol. 9, no. 11, 2008.