# EEG-based Cross-Subject Driver Drowsiness Recognition with Interpretable CNN

Jian Cui, Yisi Liu, Zirui Lan, Olga Sourina, Wolfgang Müller-Wittig

*Abstract*—In the context of electroencephalogram (EEG)-based driver drowsiness recognition, it is still a challenging task to design a calibration-free system, since there exists a significant variability of EEG signals among different subjects and recording sessions. As deep learning has received much research attention in recent years, many efforts have been made to use deep learning methods for EEG signal recognition. However, existing works mostly treat deep learning models as blackbox classifiers, while what have been learned by the models and to which extent they are affected by the noise from EEG data are still underexplored. In this paper, we develop a novel convolutional neural network that can "explain" its decision by highlighting the local areas of the input sample that contain important information for the classification. The network has a compact structure for ease of interpretation and takes advantage of separable convolutions to process the EEG signals in a spatial-temporal sequence. Results show that the model achieves an average accuracy of 78.35% on 11 subjects for leave-one-out cross-subject drowsiness recognition, which is higher than the conventional baseline methods of 53.4%-72.68% and state-of-art deep learning methods of 63.90%-65.61%. Visualization results show that the model has learned to recognize biologically explainable features from EEG signals, e.g., Alpha spindles, as strong indicators of drowsiness across different subjects. In addition, we also explore reasons behind some wrongly classified samples and how the model is affected by artifacts and noise in the data. Our work illustrates a promising direction on using interpretable deep learning models to discover meaning patterns related to different mental states from complex EEG signals.

*Keywords-EEG, convolutional neural networks, class activation map, driver drowsiness recognition, visualization technique*

## I. INTRODUCTION

Causing decrease in attention, vigilance and cognitive performance, driver's drowsiness is a leading factor of car accidents. Development of a drowsiness monitoring system to continuously watch the vigilance state of the driver and send alarm before the driver falls asleep is of high priority for safety driving and prevention of transportation accidents. Many efforts have been made to investigate monitoring driver drowsiness using electroencephalogram (EEG), which is believed to be the most practical non-invasive modality for capturing brain dynamics due to its high temporal resolution and low cost. The association between drowsiness and change of EEG signals has been extensively researched. It was found that drowsiness resulted from night driving [1] and sleep deprivation [2] can cause power increase in the Theta and Alpha frequency bands. EEG Alpha spindles, which appear as short narrowband bursts in the Alpha band, were also discovered to be indicators of driver drowsiness [3]. In addition, entropy features extracted from EEG signals were found useful to recognize drowsiness [4, 5].

Despite the progress, building a calibration-free drowsiness recognition system is still a challenging task. The difficulty lies in capturing common drowsiness-related patterns from a diversity of EEG signals with a low signal-to-noise rate. The variability of EEG signals from different subjects is attributed by many factors, such as electrode displacements, skin-electrode impedance, different head shapes and sizes, different brain activity patterns, and disturbance by task-irrelevant brain activities. Conventional methods relying on hand-crafted features are often very specific to some EEG characteristics of interest, which potentially excludes other relevant information that could be essential to drowsiness recognition. In comparison, deep learning allows end-to-end learning without need for priori feature crafting. Such models can directly learn essential characteristics from raw high-dimensional data by converting it into a cascade of representations while optimizing the parameters through back propagation. Since its initial success in many challenging image classification problems [6, 7], deep learning is receiving more and more attention, and it has also been used in the area of EEG signal recognition, e.g., sleep stages recognition [8], brain-computer interface (BCI) [9], and workload levels classification [10]. However, existing works in the area of EEG signal processing mostly treat deep learning models as blackbox classifiers, while what have been learned by the models and to which extent they are affected by the noise from EEG data are still underexplored. Without knowing these facts, the works of developing the model towards higher accuracy become a trial-and-error process.

In this paper, we propose a novel Convolutional Neural Network (CNN) for driver drowsiness recognition and discovering common drowsiness-related patterns of EEG signals across different subjects. The network has a compact structure and it uses separable convolutions to process the EEG signals in a spatial-temporal sequence. In order to allow the model to "explain" its decisions, we have designed visualization techniques specially for the model that can reveal local regions of the input signals that are 'important' for prediction. In the following part of this paper, existing works are reviewed in Section II. The methods are proposed in Section III. The performance of the proposed method is evaluated in Section IV, which is followed by discussion and future works in Section V. Conclusions are made in Section VI.

## II. RELATED WORKS

It was found that there exists a strong relationship between drowsiness and the oscillation patterns of EEG signals. For example, Akerstedt et al. [11] found a significant increase in power from the Theta (4-7.9 Hz) and Alpha (8-11.9 Hz) frequency

bands for subjects during night driving in comparation to those in day driving. Corsi-Cabrera et al. [12] found an increase in power from frequency bands of fast upper Alpha (9.77-12.45 Hz) and Beta (12.7-17.85 Hz) from subjects experienced sleep deprivation. In another experiment conducted in a driver simulator [2], the subjects were found to have higher power in the EEG frequency bands during the early sleep stage. It was summarized by Klimesch [13] that drowsiness can in general cause increase in power of the Theta and Alpha frequency bands. They further concluded that increase in lower Alpha power occurs only when subjects are struggling not to sleep, while the Alpha power will decrease when subjects fall asleep.

Conventional methods for drowsiness recognition from EEG signals mainly include two steps of feature extraction and feature classification. Feature extraction requires expertise and/or a priori knowledge for converting the data into a meaningful representation, modeling some characteristics of interest, while feature classification involves using machine learning algorithms to classify the representations into different labels. For example, Yeo et al. [14] converted EEG data into frequency domain using fast Fourier transform (FFT) and extracted four features, which are dominant frequency, average power of dominant peak, center of gravity frequency, and frequency variability, from the standard EEG bands of Delta (0.5–4 Hz), Theta (4–8 Hz), Alpha (8–13 Hz), and Beta (13–20 Hz), and used Supported Vector Machine (SVM) for classification. Hu [15] extracted entropy features from EEG signals for driver drowsiness detection, as they are widely used to quantify the complexity of the nonlinear EEG signals. Specifically, they calculated sample entropy, fuzzy entropy, approximate entropy, and spectral entropy from EEG samples combined with four classifiers, which are gradient boosting decision tree, k-nearest neighbor, support vector machine and neural network, to distinguish between EEG signals under alert and drowsy states. Considering traditional entropy features calculated from a single time scale cannot measure long distance correlation, Hu and Min [4] proposed an adaptive multi-scale entropy feature extraction algorithm to process forehead EEG data. Luo et al. [5] argued wavelet entropy has advantages over time domain entropy, and they proposed to use wavelet entropy combined with SVM for driver fatigue state monitoring.

With the emergence of deep learning, the performance and accuracy of many artificial intelligence and classification tasks in fields, such as computer vision [16] and speech recognition [17], have been vastly boosted. It has also become an emerging direction in the fields of EEG signal processing, such as brain-computer interface (BCI) [9], classification of different sleep stages [8] and workload levels [10]. In comparison to traditional EEG processing methods based on feature crafting, deep learning can directly learn from raw data and transform it into a cascade of representations with increasing abstraction. The important characteristics of raw EEG signals can therefore be maximally retained when the processes of feature extraction and classifier training are combined under the same learning framework. For example, Schirrmeister et al. [18] proposed two types of convolutional network structures to decode raw EEG data. The first model named deep ConvNet has four convolution-max-pooling blocks and a dense softmax classification layer, while the second model named shallow ConvNet has two convolutional layers performing temporal and spatial convolutions followed by a pooling and a dense layer, which mimics the filter bank common spatial patterns (FBCSP) pipeline [19]. Lawhern et al. [9] proposed a compact convolutional network called EEGNet to classify EEG signals in different BCI paradigms. The novelty of the network is the introduction of depthwise and separable convolutional layers to replace the conventional convolutional layers, which dramatically reduce the complexity of the network and can thus be trained with limited size of data. Their results show that EEGNet achieves comparably high performance to reference algorithms across different paradigms. A detailed review on current progress on deep learning for EEG signal processing can be found in [20].

As for driver drowsiness recognition, Rundo et al. [21] extracted frequency domain features using the discrete cosine transform (DCT) method [22], and used a network consisting of stacked autoencoder and Softmax layers to classify EEG signals under alert and drowsy states. Their results show that the method has a higher accuracy over conventional machine learning methods, including Support Vector Machine (SVM) and Linear Discriminant Analysis (LDA). As for end-to-end deep learning, Nissimagoudar et al. [23] proposed a simple convolutional network consisting of two convolutional layers, a max-pooling layer, a flatten and a fully connected layer to classify single-channel EEG data for Advance Driver Assistance Systems (ADAS) of automotive. Ding et al. [24] implemented a deep CNN model on mobile device to detect drowsiness from single-channel EEG signals. The model employs cascaded CNN and attention mechanism layer in the structure. Results show their model outperforms other benchmark deep learning models, including AlexNet [6] and ResNet, as well as conventional machine learning methods, such as Support Vector Machine (SVM) [25] and Linear discriminant analysis (LDA). In order to process multi-channel EEG signals, Gao et al. [26] proposed a model called spatial–temporal CNN (ESTCNN) for driver fatigue detection. The model contains three core blocks, and each block has a convolution layer, a rectified linear activation layer and a batch normalization layer. Zeng et al. [27] developed two CNN models called EEG-Conv and EEG-Conv-R, respectively. The first one is based on the traditional CNN and second one combines CNN with deep residual learning. They found both models outperform the LSTM- and SVM-based classifiers, while EEG-Conv-R converges more quickly. Hajinoroozi et al. [28] designed a network called channel-wise convolutional neural network (CCNN) and a variation of the model by replacing the convolutional filters with Restricted Boltzmann Machine. They found both models had improved performance over conventional methods and deep learning methods.

Although existing works show promising results on using deep learning to recognize driver drowsiness from EEG signals, there is still little insight on what characteristics of the EEG data have been learned by the models to distinguish the signals between different mental states and to which extent the deep learning models are affected by noise from the data. In fact, deriving insights on what has been learned by the deep learning networks has become an important procedure of model validation, since it can not only ensure the classification is driven by relevant features rather than noise or artifacts in the data but also potentially

discover interesting neurophysiological phenomena that explain the model decisions. In this connection, we design a compact CNN model that can "explain" its decisions for the purpose of learning and understanding common EEG patterns related to different mental states in the task of driver drowsiness recognition.

## III. MATERIALS AND METHODS

### A. Data preparation

A public EEG dataset [29] was used in the study. The dataset was collected from 27 subjects (aged between 22–28), who were students or staff from the National Chiao Tung University. The EEG signals were sampled at 500 Hz with 30 electrodes and processed with 1-50 Hz bandpass filters and artifact rejection. The pre-processed version of the dataset available from [30] was used in this study. We further down-sampled the original data to 128 Hz and extracted the EEG samples of 3-second length prior to the car deviation events for each trail. Each sample has a dimension of 30 (channel) × 384 (sample points). We followed methods described in [31] to select and label the EEG samples. Specifically, the local reaction time (RT), which is the time taken by the subject to respond to the car drift event, and the global-RT, which is the average of RTs within a 90-second window before the car drift event, were calculated for each sample. The baseline 'alert-RT' for each session was defined as the 5th percentile of the local RTs. Samples with both local and global-RT shorter than 1.5 times alert-RT were labeled as alert state, while samples with both local and global RT longer than 2.5 times alert-RT were labeled as drowsy state. We discarded sessions with less than 50 samples of either class. If there were multiple sessions of the same subject, the session with the most balanced class distribution for the subject was used. Samples from each session were further balanced by choosing the most representative ones from the majority class according to their local RTs. In this way, we finally got 2022 samples in total from 11 different subjects. The number of samples for each subject/session is shown in Table 1.

**Table 1**. Number of extracted samples from each eligible subject

| Subject ID | Sample Number | |
|:---:|:---:|:---:|
| | Alert | Drowsiness |
| 1 | 94 | 94 |
| 2 | 66 | 66 |
| 3 | 75 | 75 |
| 4 | 74 | 74 |
| 5 | 112 | 112 |
| 6 | 83 | 83 |
| 7 | 51 | 51 |
| 8 | 132 | 132 |
| 9 | 157 | 157 |
| 10 | 54 | 54 |
| 11 | 113 | 113 |
| Total | 1011 | 1011 |

### B. Network design

The EEG signals can be viewed as a mixture of cortical source signals generated from different areas of the brain. However, the recorded data are inevitably contaminated by artifacts caused by different electrical activities, e.g., cardiac, eye movement and muscle tension, as well as noise generated from the equipment. Considering learning directly from the noisy and redundant EEG data usually leads to unsatisfactory recognition results, spatial filtering techniques [32] were proposed to improve the data quality by extracting a set of new signals from the raw multi-channel recordings of the EEG data with minimal contamination and redundancy. Specifically, suppose the EEG signals recorded from $m$ electrodes are $\{x_i\}_{i=1,2\ldots m}$. $N_1$ new signals $\{s_j\}_{j=1,2\ldots N_1}$ can be obtained from linear combination of the original $m$ signals.

$$s_j = \sum_{i=1}^{m} w_{i,j} x_i + b_j \tag{1}$$

In Equation (1), the weights $\{w_{i,j}\}$ are a set of spatial filters, which can be calculated based on various independent evaluation criteria such as distance measure, information measure, dependency measure, and consistency measure [32]. For example, Independent Component Analysis (ICA) [33] is one of the most well-known methods that finds the weights by solving an equation based on the statistical independency hypothesis of the source signals. Other methods are Common Spatial Pattern (CSP) [34], Minimal Energy Combination (of noise), Maximum Contrast Combination (MCC) [35], Canonical Correlation Analysis (CCA) [36], and so forth. Since the obtained new set of signals $\{s_j\}$ are expected to contain least noise and redundancy, a set of features can be thus extracted from each new signal $S_j$ to be used for classification.

$$[feature_{j,1} \; feature_{j,2} \ldots, feature_{j,k}] = f_j(S_j) \tag{2}$$

Based on the observations above, we consider designing a compact CNN model that processes EEG data in a similar spatial-temporal sequence. The network consists of seven layers and its structure is shown in Figure 1. The processing steps described in Equation (1) and (2) are implemented in the model with pointwise and depthwise convolutions, respectively. The spatial filters $\{w_{i,j}\}$ and parameters of $\{f_j\}$ are treated as trainable network parameters and they are updated simultaneously in the process of optimizing the network. Specifically, in the first layer we use $N_1$ 1D pointwise convolutional nodes to generate $N_1$ new signals. For a given input sample $X_{m \times n}$, where $m=30$ and $n=384$. The outputs from the first layer is

$$h_{i,j}^I = \sum_{p=1}^{m} w_{i,p}^I x_{p,j} + b_i^I, \tag{3}$$

where $i = 1, 2, 3, …, N_1$ and $x_{p,j}$ is a sampling point of an EEG sample from $X_{m \times n}$. $w_{i,p}^I$ and $b_i^I$ are the weight and bias in the first layer, respectively. The Roman numerals superscripts of the outputs and network parameters indicate the number of layer that they belong to. We set $N_1=16$, which is around half of the input channels, in order to reduce redundancy and encourage convergence of the network.

In the second layer, depthwise convolutions are used to extracted features from the $N_1$ extracted signals. Specifically, each new signal $h_i^I$ is convoluted with two nodes in the second layer. Suppose the length of the kernel is $l$, the output of the layer is

$$h_{i,j}^{II} = \begin{cases} \sum_{r=1}^{l} h_{\frac{i+1}{2},j+r-1}^I w_{i,r}^{II}, for\ i\ is\ odd. \\ \sum_{r=1}^{l} h_{\frac{i}{2},j+r-1}^I w_{i,r}^{II}, for\ i\ is\ even. \end{cases} \tag{4}$$

In Equation (4), the length of a kernel is set as $l = 64$, which is half of the sampling rate (128 Hz). The size of the output $h_{i,j}^{II}$ is $(2N_1, m-l+1)$, which is (32, 321).

The 3$^{th}$ and 4$^{th}$ layers are activation and batch normalization layers.

$$h_{i,j}^{III} = ReLU(h_{i,j}^{II}) \tag{5}$$

$$h_{i,j}^{IV} = BatchNorm(h_{i,j}^{III}) \tag{6}$$

Global Average Pooling (GAP) [37] is used in the 5$^{th}$ layer. In comparison to the widely used fully connected layer, the GAP layer dramatically reduces parameters and can thus effectively prevent over-fitting.

$$h_i^V = (\sum_{j=1}^{m-l+1} h_{i,j}^{IV})/(m - l + 1) \tag{7}$$

The model ends with a dense layer and a Softmax activation layer.

$$h_c^{VI} = \sum_{i=1}^{2N_1} w_{i,c}^{VI} h_i^V + b_{i,c}^{VI} \tag{8}$$

$$h_c^{VII} = Softmax(h_c^{VI}) \tag{9}$$

In equation (8) and (9), $c = 0$ or 1, which represents the alert or drowsy state, respectively.

In the network design, separable convolution consisting of a pointwise convolution (1$^{st}$ layer) and depthwise convolution (2$^{nd}$ layer) is used to implement the 2-step spatial-temporal sequence of EEG signal processing. Actually, separable convolution has been previously used in deep CNN models [38] to replace standard convolutions for the purpose of reducing parameters and
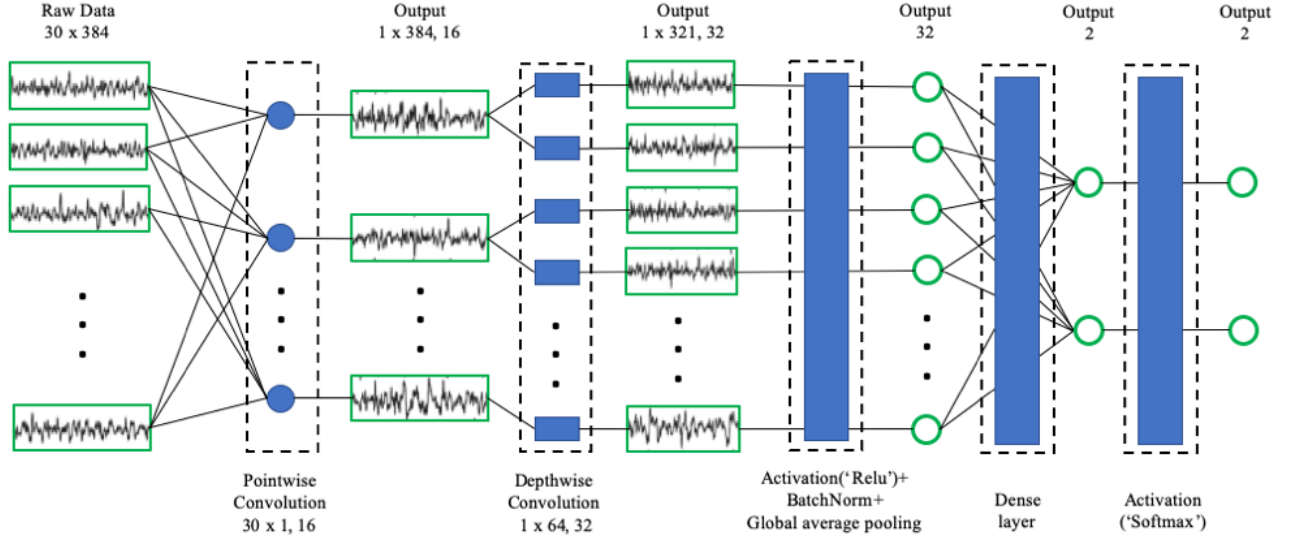
Figure 1. The architecture of the proposed model. The shapes with green outlines indicate the dimension changes of the EEG samples in the intermediate layers. The blue shapes inside the dashed borders represent kernels or layers of the network.

encourage convergence of the network. In our model, the structure is used with the expectation to extract essential features from EEG signals that can best distinguish between alert and drowsy EEG signals. The spatial-temporal processing pipeline is similar to the method proposed by Lin et al. [39], where ICA is applied to demix raw EEG data and band power features are extracted from each of the demixed signals separately, which are used to estimate driver's drowsiness. By comparison, we have incorporated the processes in the CNN model structure allowing the parameters to be optimized towards a high recognition accuracy.

*C. Visualization technique*

Deriving insights on what characteristics of EEG signals has been learned by the deep learning networks has become an important procedure of model validation, since it can not only ensure the classification is driven by relevant features rather than noise or artifacts in the data but also potentially discover interesting neurophysiological phenomena that explain the model decisions. Existing works using deep learning for EEG signal classification have attempted to interpret the model by visualizing kernel weights, summarizing averaged output from hidden unit activations, calculating single-trial feature relevance [9], reconstruction of the sample that leads to maximal activations [40], and so forth. Although these techniques allow understanding of what global patterns have been learned from the massive data, they cannot explain what specific characteristics of each EEG sample have been found relevant to different mental states and to which extent the models are affected by noise from the data.

By comparison, the Class Activation Map (CAM) [41] method is a powerful technique that can localize the discriminative regions of a sample for a CNN model trained to solve a classification task. Specifically, for each input sample a heatmap is generated from the activations after the last convolutional layer. The map is then interpolated to size of the input sample and it reveals to which extent the local regions of the input sample contribute to the classification. The method can potentially allow us to improve the model performance by revealing whether biologically explainable features have been identified for classification. However, the CAM method was originally designed for deep CNN networks with only standard convolutional layers (e.g., Alexnet [6]) for classification of image data, while it cannot be directly used for the proposed model, as well as other CNN models involving convolutions in the spatial dimensional [26] or having structures more than basic convolutions blocks [9]. Therefore, we start from the CAM method and consider designing a visualization technique specially for the proposed model.

Suppose a given input EEG sample $X_{m \times n}$ is classified with label $c$, where $c$ is either 0 or 1 representing the alert or drowsy state, respectively. The objective is to find the heatmap $S_{m \times n}^{c}$ for $X_{m \times n}$ that can reveal important regions for the prediction by the network. Suppose an input signal $X_{m \times n}$ generates activation $h_{c}^{VI}$ in the 6th layer of the network. From equation (7) and (8), we have

$$h_{c}^{VI} = \sum_{i=1}^{2N_1} w_{i,c}^{VI} h_{i}^{V} = \sum_{i=1}^{2N_1} \sum_{j=1}^{m-l+1} w_{i,c}^{VI} h_{i,j}^{IV} = \sum_{i=1}^{2N_1} \sum_{j=1}^{m-l+1} M_{i,j}^{c} \ , \tag{10}$$

where

$$M_{i,j}^c = w_{i,c}^{VI} h_{i,j}^{IV} \qquad (11)$$

$M_{i,j}^c$ is the activation map of class $c$ for the sample $X_{m \times n}$, as defined in [41]. In Equation (10), we neglect the constants $b_{i,c}^{VI}$ and $(m - l + 1)$ for simplicity. As it can be seen in Equation (10), $M_{i,j}^c$ can be viewed as the distribution of the final activation $h_c^{VI}$ for class $c$ in a map of size $2N_1 \times (m - l + 1)$. The original CAM method finds the heatmap by upsampling $M_{i,j}^c$ until it has the same size as the input sample. However, the method cannot be directly used for our model since the channels of the input signal are mixed by pointwise convolutions in the first layer, which makes the first (channel) dimension of $M_{i,j}^c$ misaligned with the first (channel) dimension of the input signal. Inspired by the CNN-Fixation method [42], we consider an alternative way of tracing only a small portion of the positions in the activation map $M_{i,j}^c$ that contribute most to the class activation $h_c^{VI}$ rather than the whole activation map, back to the their major corresponding areas in the input signal. Specifically, we rank the values of $M_{i,j}^c$ in a descending order. Suppose the locations of the first $N$ elements in $M_{i,j}^c$ are $(i_0, j_0)$, $(i_1, j_1)$, ..., $(i_N, j_N)$, where $1 \leq i_k \leq 2N_1 (2N_1 = 32)$ and $1 \leq j_k \leq m - l + 1$ $(m - l + 1 = 321)$. The objective is to trace each of these discriminative locations for class $c$ in $M_{i,j}^c$ throughout the network to the center of areas in the input sample that contribute most to these high activations. Suppose we can find $N$ corresponding discriminative locations in the input samples and they are $(p_0, q_0)$, $(p_1, q_1)$, ..., $(p_N, q_N)$, where $1 \leq p_k \leq m$ $(m = 30)$ and $1 \leq q_k \leq n$ $(n = 384)$. The final heatmap for sample $X_{m \times n}$ can be obtained by combining all the class discriminative points in the input sample with the Gaussian function.

$$S_{p,q}^c = \frac{1}{\sigma\sqrt{2\pi}} \sum_{k \; for \; p_k = p} e^{\left(-\frac{1}{2}\frac{(q - q_k)^2}{\sigma^2}\right)} \qquad (12)$$

In equation (12), $\sigma$ is a constant that decides radius of the influential area of each discriminative point in the input signal. $S_{p,q}^c$ is further normalized in the range (-1, 1) for visualization.

Finally, we consider how to trace the set of discriminative locations $(i_0, j_0)$, $(i_1, j_1)$, ..., $(i_N, j_N)$ in $M_{i,j}^c$ back to the input sample through the four layers (layer 1-4) of the proposed network. It is easy to notice that the discriminative locations are unchanged after the 3$^{rd}$ and 4$^{th}$ layers of the network, since the activation and batch normalization layers only perform element-wise operations that will not affect the topology of the data. Therefore, the only task left for us is to consider how to trace the discriminative locations through the depthwise and pointwise convolutional layers, which is not discussed in the original CNN-Fixation method [42]. Suppose the input sample $X_{m \times n}$ generates activation $h_{i_k,j_k}^{II}$ after the 2$^{nd}$ layer of the network at the discriminative location $(i_k, j_k)$. From Equation (3) and (4), we have

$$h_{i_k,j_k}^{II} = \sum_{r=1}^{l} h_{\frac{i_k+1}{2},j_k+r-1}^{I} w_{i_k,r}^{II}$$

$$= \sum_{r=1}^{l} w_{i_k,r}^{II} \sum_{p=1}^{m} w_{\frac{i_k+1}{2},p}^{I} x_{p,j_k+r-1}$$

$$= \sum_{p=1}^{m} w_{\frac{i_k+1}{2},p}^{I} \sum_{r=1}^{l} w_{i_k,r}^{II} x_{p,j_k+r-1} , \qquad (13a)$$

when $i_k$ is odd, and similarly

$$h_{i_k,j_k}^{II} = \sum_{p=1}^{m} w_{\frac{i_k}{2},p}^{I} \sum_{r=1}^{l} w_{i_k,r}^{II} x_{p,j_k+r-1} , \qquad (13b)$$

when $i_k$ is even. In Equation (13a) and (13b), we ignore $b_i^I$ for simplicity of expression. From Equations (13a,b), we can observe that $h_{i_k,j_k}^{II}$ is generated from an episode of the input signals at the local area from the time point $j_k$ to $j_k + l - 1$. Actually, it is the weighted sum of the convoluted signals $\sum_{r=1}^{l} w_{i_k,r}^{II} x_{p,j_k+r-1}$ of all the $m$ channel, and the weight assigned to the channel $p$ is

$w^I_{\frac{i_k+1}{2},p}$ or $w^I_{\frac{i_k}{2},p}$. Therefore, the discriminative locations $(i_k, j_k)$ in $M^c_{i,j}$ can be traced back to the center $(p_k, q_k)$ of the strongest contributing episode in the input signal, where

$$p_k = argmax_p \left( w^I_{\frac{i_k+1}{2},p} \sum_{r=1}^{l} w^{II}_{i_k,r} x_{p,j_k+r-1} \right) \tag{14a}$$

when $i_k$ is odd,

$$p_k = argmax_p \left( w^I_{\frac{i_k}{2},p} \sum_{r=1}^{l} w^{II}_{i_k,r} x_{p,j_k+r-1} \right) \tag{14b}$$

when $i_k$ is even, and

$$q_k = j_k + (l-1)/2 \tag{15}$$

We set $\sigma = l/2 = 32$ for Equation (12), so that the discriminative location in the input signal will highlight the whole episode of strongest contributing signal. We trace the top 100 (*N*=100) discriminative locations in class activation map $M^c_{i,j}$, which accounts for around 1% of all entries of $M^c_{i,j}$.

### D.  Methods for comparison

In this part, we compare the performance of the proposed model with both state-of-art deep learning and conventional baseline methods. In order to understand how each part of the model influences its performance, we also compare the model with its variations, where a single component of the model is replaced or deleted in the structure.

### 1)  Deep learning methods

The deep learning model we used for comparison is the benchmark CNN model for EEG signal classification—EEGNet proposed by Lawhern et. al. [9]. Inspired by the filter bank common spatial patterns (FBCSP) algorithm [43], EEGNet uses a standard convolutional layer to filter the raw signals, which is followed by a depthwise convolutional layer to extract spatial features from the filtered signals. The model was tested on several Brain Computer Interface (BCI) datasets and achieved higher accuracies over conventional methods. The model has also been tested on cross-subject driver drowsiness recognition in a preliminary study conducted by Liu et al. [44], where two configurations of EEGNet – EEGNet-4,2 and EEGNet-8,2 were used as baseline deep learning methods for comparison.

### 2)  Conventional baseline methods

Conventional methods for EEG signal classification mainly involve the stages of feature extraction and feature classification. In order to have a comprehensive understanding on the performance of different conventional methods on the dataset, we implemented five baseline methods for feature extraction and tested them on eight different classifiers for comparison.

EEG band power features have been regarded as golden standard for EEG signal classification. For driver drowsiness recognition, many works [2, 12, 45] have found a strong relationship between drowsiness and band power features of EEG signals. Therefore, the first three baseline methods use different forms of band power features, which are relative band power features, log of band power features [31], and the ratio of band power features [46]. The third baseline method uses the wavelet entropy features [47], while the forth baseline method uses a combination of four entropies features [4], which are sample entropy, fuzzy entropy, approximate entropy and spectral entropy. Details of these methods are illustrated below.

**RelativePower**: We use relative band power features as the first baseline method for comparison. Specifically, power features from four frequencies bands of Delta (1–4 Hz), Theta (4–8 Hz), Alpha (8–12 Hz) and Beta (12–30Hz) are extracted from each EEG channel. Considering the absolute values of the band powers vary significantly across different samples, relative power of the four frequency bands from each EEG channel are calculated.

**LogPower**: The second baseline method is slightly different from the first one—natural log of the band power instead of relative power is calculated. The method was proposed by Pal et al. [31] based on the observation of a strong linear correlation between log power features of EEG and subject's driving performance.

**PowerRatio**: Jap et al. [46] found four band power ratios (i) $(\theta + \alpha)/\beta$ , (ii) $\alpha/\beta$, (iii) $(\theta + \alpha)/(\alpha + \beta)$ and (iv) $\theta/\beta$ were good indicators of driver drowsiness. Therefore, for the third baseline method, we calculated the four band power ratio features and use them as representations of the EEG sample signals.

**WaveletEntropy**: We implemented the method proposed by Wang et al. [47], where wavelet entropy features from EEG signals are used to recognize driver drowsiness. Specifically, the Mexican Hat Wavelet is used in our implementation, and wavelet coefficients on the wavelet scales of 0.5, 1, 2, 4, 8, 16, 32 (corresponding to frequencies of 64 Hz, 32 Hz, 16 Hz, 8 Hz, 4 Hz, 2 Hz, 1 Hz) are extracted from each EEG channel. The wavelet entropy feature for each EEG channel is calculated by applying the Shannon function on the normalized wavelet coefficients.

**FourEntropies**: Hu et al. [4] proposed to use four types of entropies, which are sample entropy, fuzzy entropy, approximate entropy and spectral entropy for driver fatigue recognition. Following the descriptions in the paper, we calculated the approximate entropy and spectral entropy using the methods proposed by Song et al. [48] and the fuzzy entropy by the method proposed by Xiang et al. [49], and set the parameters $m$ and $r$ involved in the calculation as $m = 2$ and $r = 0.2 * SD$. We set the width of the exponential function $n$ as $n=2$ for extracting fuzzy entropy features. Finally, we normalized each feature dimension for each subject, as indicated in that paper.

**Classifiers**: Different classifiers were implemented, which include Decision Tree (DT), Random Forest (RF), k-nearest neighbors (KNeighbors), Gaussian Naive Bayes (GNB), Logistic Regression (LR), Linear Discriminant Analysis (LDA), Quadratic Discriminant Analysis (QDA), and SVM.

*3)  Variations of the model for comparison*

The proposed model is compared with its variations in order to understand how each part of the model influences its overall performance. Firstly, we want to evaluate the benefits of using the first layer (containing a set of spatial filters), as well as the advantage of using the separable convolution ($1^{st}$ and $2^{nd}$ layers) over standard convolution. Therefore, in the first variation we remove the first layer of the network so that each input channel of EEG signal is directly convoluted with two kernels in the depthwise convolutional layer. In the second variation, the first two layers of the original network are replaced with a single layer of standard convolution containing 32 convolutional kernels with dimensions of $30 \times 64$.

Secondly, we consider the variations where an additional activation or batch normalization layer is added between pointwise convolutional layer ($1^{st}$) and the depthwise convolutional layer ($2^{nd}$). State-of-art CNN models commonly use an activation layer after the convolutional layers to add non-linearity transformation on the data or a batch normalization layer to remove the internal covariate shifts [26]. However, Chollet [38] found that adding a non-linear activation layer between the pointwise and depthwise convolutions will deteriorate the performance of their deep learning model designed for image classification. In order to understand whether an intermediate layer is necessary, we consider three variations and each of them has an additional batch normalization layer, an ELU activation layer, or a ReLU activation layer, respectively, after the first layer of the original network.

Thirdly, we want to evaluate whether the ReLU activation layer (the $3^{rd}$ layer) and batch normalization layer (the $4^{th}$ layer) are optimal in design for the network. Farahat et al. [50] found that their CNN model designed for decoding EEG signals in a covert attention task achieved best performance with Tanh and ELU activation layers, while using ReLU activation layer instead leads to decrease of accuracy. They also found the effect of batch normalization is mainly negligible. ELU activation is also preferred over ReLU activation in some benchmark CNN models [9, 18] for EEG signal classification. Therefore, in order to figure out which kind of activation layer is best for the proposed model, we consider to compare with variations where the ReLU activation is replaced with ELU or Tanh activation. We also test the variations without the activation or the batch normalization layer.

Lastly, we want to evaluate the effect of the GAP layer of the network. We compare the GAP layer with the commonly used standard average pooling layer. We test three variations with pooling sizes of 20, 40 and 80. Since using average pooling layer will make number of output nodes after the layer increase to 16, 8 and 4 times of the original number, we set dropout rates (probability of the node to be dropped out) after the pooling layer as 0.9375, 0.875, and 0.75 for the three variations, respectively. The names of all the variations are listed below:

**StandardConv**: the first two layers of the original network are replaced with a single layer of standard convolution containing 32 convolutional kernels with dimension of 30 x 64.
**NoSpatialFilters**: the pointwise convolutional layer ($1^{st}$ layer) of the network is removed.
**AddBatchNorm**: a batch normalization layer is added between the first two layers of the network.
**AddELU**: an ELU activation layer is added between the first two layers of the network.
**AddReLU**: a ReLU activation layer is added between the first two layers of the network.
**ELU**: the ReLU activation of the network ($3^{rd}$ layer) is replaced with an ELU activation layer.
**Tanh**: the ReLU activation of the network ($3^{rd}$ layer) is replaced with a Tanh activation layer.
**NoActiv**: the ReLU activation of the network ($3^{rd}$ layer) is removed.
**NoBatchNorm**: the batch normalization layer ($4^{th}$ layer) of the network is removed.
**AvePool20**: the GAP layer ($5^{th}$ layer) of the original network is replaced with a standard average pooling layer with the pooling size of 20. The dropout rate after this layer is set as 0.9375.

**AvePool40**: the GAP layer ($5^{th}$ layer) of the original network is replaced with a standard average pooling layer with the pooling size of 40. The dropout rate after this layer is set as 0.875.

**AvePool80**: the GAP layer ($5^{th}$ layer) of the original network is replaced with a standard average pooling layer with the pooling size of 80. The dropout rate after this layer is set as 0.75.

### E. Implementation details

The comparison was conducted on an Alienware Desktop with 64-bit Windows 10 operation system powered by Intel(R) Core(TM) i7-6700 CPU and an NVIDIA GeForce GTX 1080 graphics card. The codes were implemented and tested on the platform of Python 3.6.6. The proposed model and its variations were implemented with the Pytorch Library. The EEGNet models were downloaded from [51] and run with the Keras API of TensorFlow. As for the conventional methods, band power features were extracted using the Welch method from the SciPy library [52]. The classifiers were implemented with the sklearn library [53] and the default parameters were used.

We conducted leave-one-subject-out tests on the methods. Specifically, the EEG data from one subject was used for testing, while data from all the other subjects were used for training the classifiers. The process was iterated until every subject served once as the test subject. For training of the neural network models, we set batch size as 50 and used Adam method [54] with default parameters ($\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$) for optimization. We trained the network models from 1 to 50 epochs and evaluated the results for each epoch. Considering neural networks are stochastic, we repeated the process for 10 times. We randomized the network parameters in the beginning of each iteration. In this way, 10 (times) x 11 (subjects) = 110 folds were created for each epoch.

## IV. EVALUATION ON THE PROPOSED METHOD

### A. Model comparison results

#### 1) Comparison with baseline methods

The accuracies of the proposed model and two EEGNet models against training epochs from 1 to 50 are shown in Figure 2. As it can be seen in the figure, the proposed model has an overall better performance than the other two models. It reaches the peak accuracy of 78.35% after 11 epochs, and then it drops a little but still stabilizes at above 76% in the rest of the first 50 epochs. The accuracy trends of EEGNet-4,2 and EEGNet-8,2 follow the same pattern—both models reach their highest accuracies, which are 65.61% and 63.90%, respectively, after 3 epochs. However, the accuracies drop significantly to as low as around 54% from the $4^{th}$ epoch to the $13^{th}$ epoch, after which the accuracies rise slowly and stabilize at only around 59%. We compared the mean accuracies between the proposed model and the EEGNet models at the $3^{rd}$, $11^{th}$, and $20^{th}$ epoch. Paired t-tests show that the mean accuracies of the proposed model are significantly higher than those of both the EEGNet models at the measured epochs ($p \approx 0$).
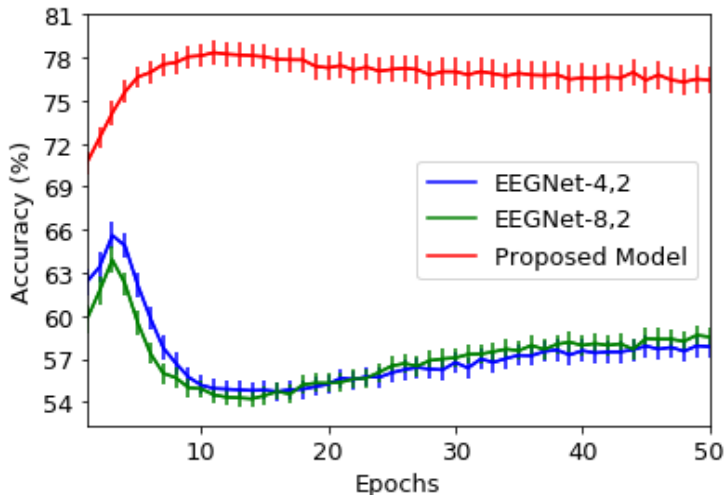


Figure 2. Cross-subject classfication accuracies (%) of EEGNet-4,2, EEGNet-8,2 and the proposed model against training epochs from 1 to 50. The standard errors and averaged accuracies over 10 iterations for 11 subjects of each model are shown.

**Table 1.** The mean cross-subject classification accuracies (%) of the five baseline methods combined with different classifiers.

| | RelativePower | LogPower | PowerRatio | WaveletEntropy | FourEntropies |
|---|---|---|---|---|---|
| DT | 60.61 | 64.30 | 60.27 | 53.40 | 58.16 |
| RF | 64.76 | 69.54 | 63.39 | 56.69 | 61.82 |
| KNeighbors | 62.66 | 71.77 | 61.62 | 57.44 | 61.95 |
| GNB | 64.76 | **72.68** | 58.75 | 56.34 | 62.96 |
| LR | 68.58 | 70.24 | 63.17 | **60.40** | 60.62 |
| LDA | 66.29 | 70.44 | 64.19 | 59.71 | 60.98 |
| QDA | 65.00 | 61.62 | 59.19 | 59.37 | 57.40 |
| SVM | **68.64** | 71.95 | **64.24** | 60.18 | **66.49** |
| Mean | 65.16 | 69.07 | 61.85 | 57.94 | 61.30 |

The accuracies of the baseline methods with different classifiers are shown in Table 1. We can see from the table that the mean classification accuracies obtained with different conventional methods range from 53.4%-72.68%. The highest mean accuracy of 72.68% is achieved by RelativePower+SVM. The best accuracies for the baseline methods of RelativePower, PowerRatio and FourEntropies, are obtained with the SVM classifier, which are 65.16%, 61.85% and 61.30%, respectively. The best accuracy for the WaveletEntropy method is achieved with the LR classifier, which is 60.40%. We also notice that the band power related baseline methods of RelativePower, LogPower and PowerRatio have an overall better accuracy over the other two baseline methods using entropy features.

Next, we compare the accuracies for each subject between the proposed model with the three baseline methods of RelativePower, LogPower and FourEntropies with their corresponding best classifiers, and the results are shown in Table 2. The mean accuracies of the proposed model for each subject are obtained by averaging over 10 repetitions after 11 training epochs. As it can be seen in Table 2, the proposed model has the highest mean accuracy of 78.35%, which is 5.67% higher than the best baseline method LogPower+GNB with mean accuracy of 72.68%. Paired t-tests show that the mean accuracy of the proposed model is significantly higher than that of RelativePower+SVM ( $p < 0.05$ ), LogPower+GNB ( $p < 0.05$ ) and FourEntropies+SVM ($p < 0.05$). It can also be observed that most of the highest individual accuracies (8 out of 11) are achieved by the proposed model.

**Table 2**. Comparison of the mean cross-subject accuracies (%) between the proposed model and three baseline conventional methods—RelativePower+SVM, LogPower+GNB, and FourEntropies+SVM. The accuracies obtained for each subject and overall mean accuracies are shown in the table.

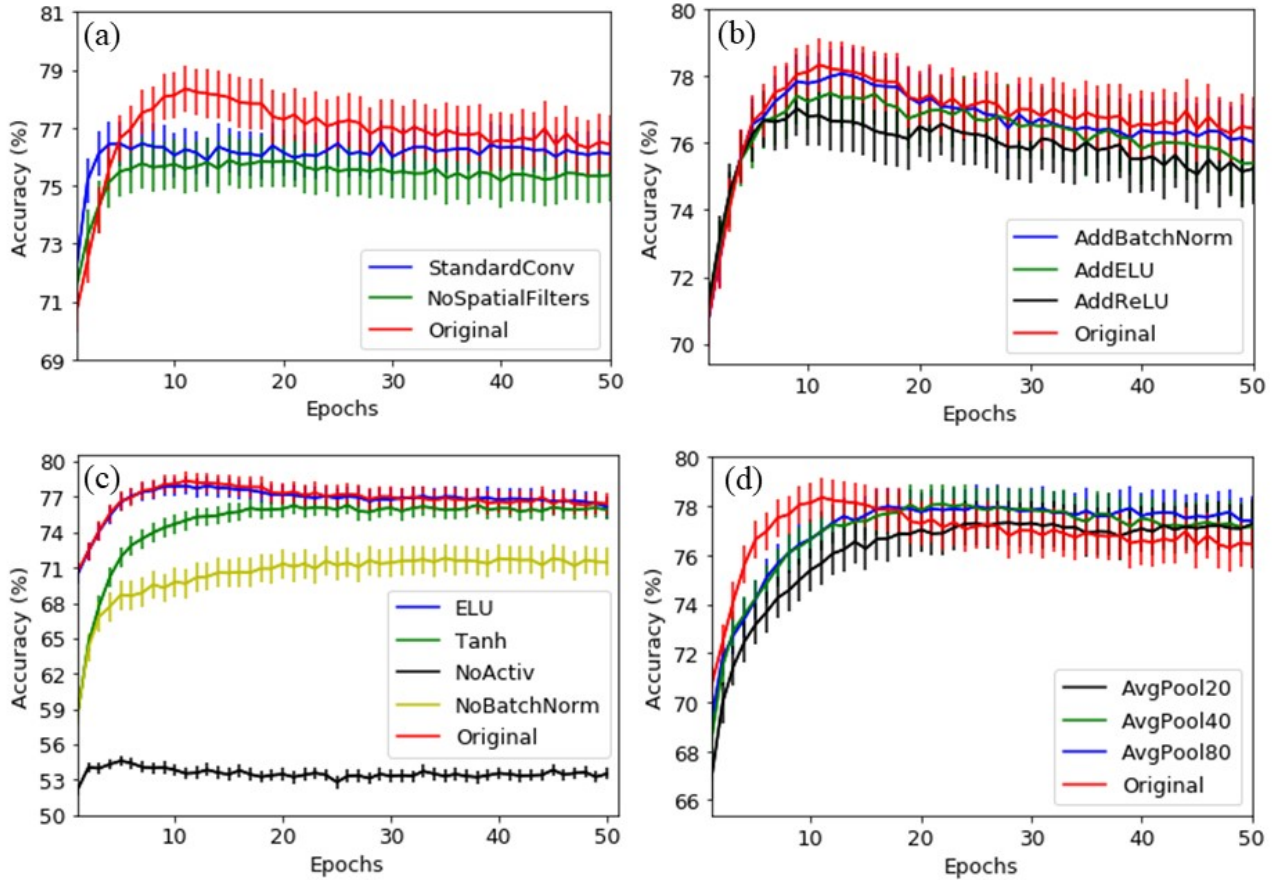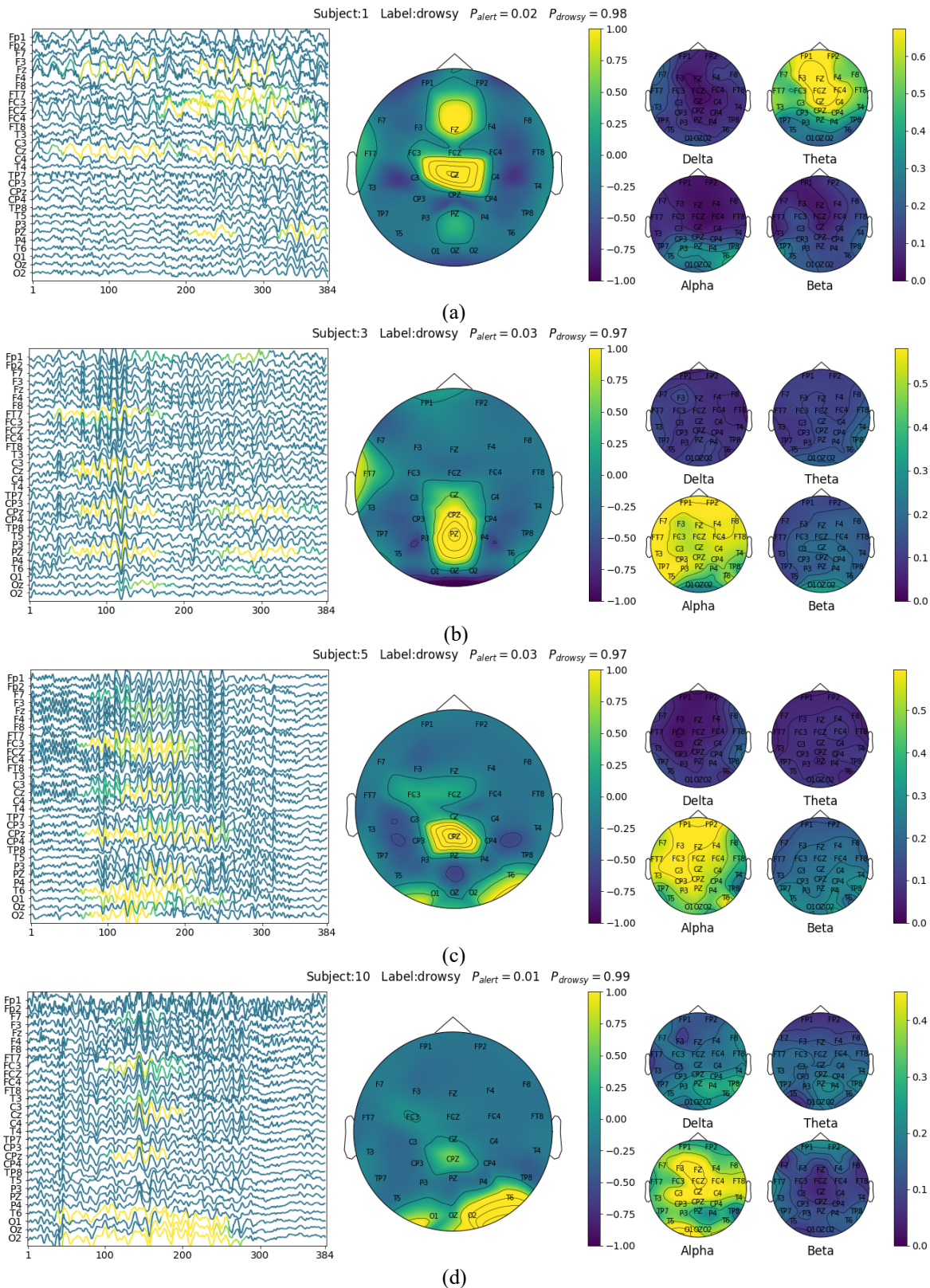| Subject ID | Relative Power + SVM | Log Power +GNB | Four Entropies +SVM | The proposed model |
|---|---|---|---|---|
| 1 | 63.30 | 78.72 | 78.19 | **85.00** |
| 2 | 53.03 | **78.79** | 58.33 | 67.65 |
| 3 | 56.67 | 65.33 | 74.00 | **81.80** |
| 4 | 58.11 | 77.70 | 39.19 | **78.99** |
| 5 | 73.21 | 75.45 | 62.95 | **88.35** |
| 6 | 83.13 | 76.51 | 67.47 | **83.92** |
| 7 | 66.67 | 56.86 | 56.86 | **67.06** |
| 8 | 73.86 | 62.88 | 69.70 | **79.05** |
| 9 | 81.85 | 87.26 | 76.11 | **89.17** |
| 10 | **82.41** | 74.07 | 77.78 | 71.02 |
| 11 | 62.83 | 65.93 | **70.80** | 69.82 |
| Average | 68.64 | 72.68 | 66.49 | **78.35** |

Figure 3. Cross-subject classification accuracies (%) of the proposed model and its variations against training epochs from 1 to 50. The standard errors and averaged accuracies over 10 iterations for 11 subjects of each model are shown.

### 2) Investigation on impact of model components

Having obtained the results from the previous section, we proceed to investigate how each part of the model impacts its overall performance. We first investigate the impact of the first two layers—the pointwise and depthwise convolutional layers. The first layer contains a set of spatial filters to extract a set of new signals from the input EEG data. When the first layer is deleted from the original network, the mean accuracy drops noticeably in comparison to the original model from the 3$^{rd}$ epoch and afterwards, as it can be seen in Figure 3(a). The results validate the usefulness of the first layer that performs spatial filtering on the multi-channel EEG signals. For the second variation where the first two layers of the network are replaced with a standard convolutional layer, the model converges at an accuracy (around 76%) lower than that of the original one, as it can be seen in Figure 3(a). The results indicate that it is optimal to use separable convolutions with no intermediate layers for the proposed model, while adding a batch normalization or an activation in-between will deteriorate the model performance, as it can be seen in Figure 3(b).

The importance of the activation layer (3$^{rd}$ layer) and the batch normalization layer (4$^{th}$ layer) after the separable convolution is indicated by the results shown in Figure 3(c)—removing either of them will significantly deteriorate the accuracies of model. It can also be observed in Figure 3(c) that the model performance is not much affected when the ReLU activation is replaced with ELU activation. However, the performance is negatively affected when the ReLU activation is replaced with Tanh activation. The obtained results on the proposed model contradict with the findings by Farahat [50] et al., where Tanh activation leads to better performance over ReLU activation for their CNN model. As it can be seen in Figure 3(d), the GAP layer has advantage over standard average pooling layers by allowing the model to converge faster to a higher accuracy. Actually, the GAP layer is in natural robust to spatial translations of the input, which makes it able to sensitively detect microstructures (e.g., EEG spindles) nested in the EEG signals that could be indicators of drowsiness.

**Figure 4.** Visualization of learned patterns on selected drowsy EEG samples that are correctly classified by the network with high likelihood. The subject ID, sample label, likelihood output by the model for alert and drowsy labels are shown on top of each sub-figures. In the left part of each sub-figure, the contributing regions to classification are highlighted by the heatmap overlaid on the input EEG signal, which is obtained

with the visualization technique described in Section III.C. The topologic heatmap in the middle of each sub-figure is obtained by averaging the heatmap over each EEG channel. It summarizes to which extent each channel contributes to the final classification. The relative powers of Delta, Theta, Alpha and Beta frequency bands for each EEG channel of the input signal is shown in the right part of each sub-figure.
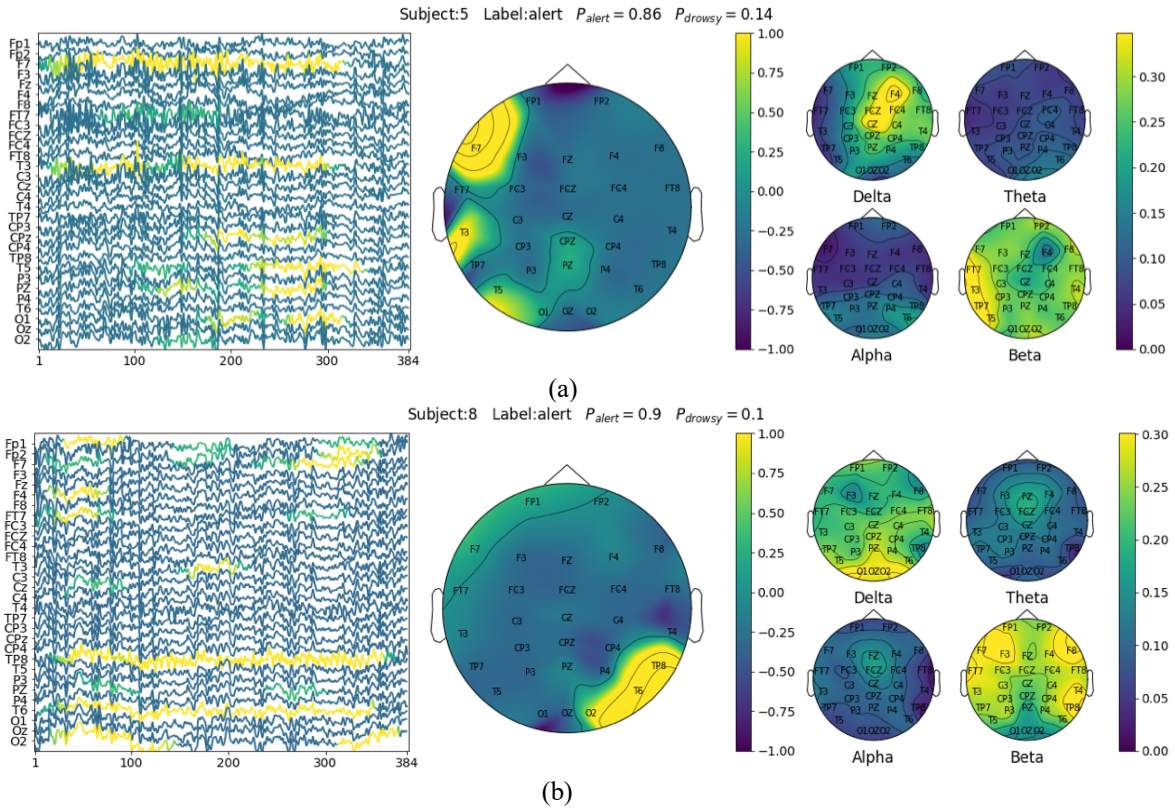
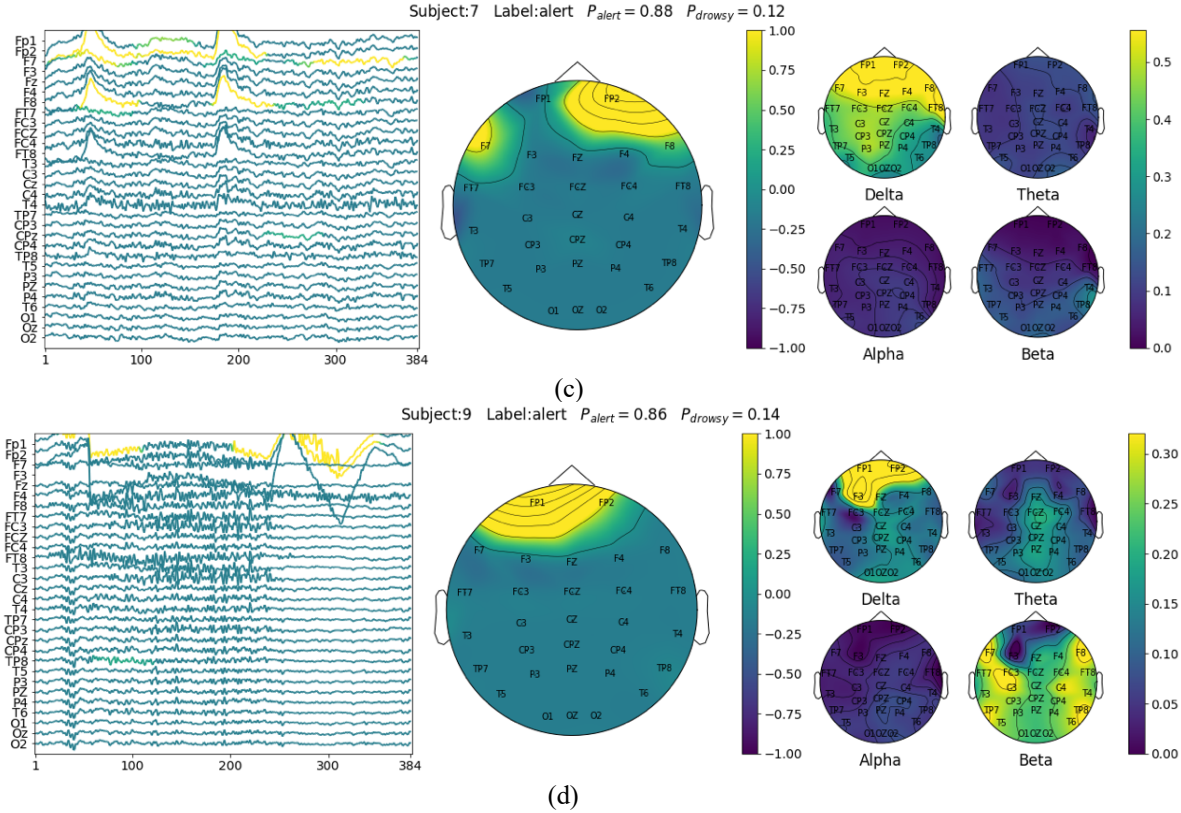### B. Visualization on the learned characteristics from EEG signals

Deriving insights into what the model has learned from the data is an important procedure of model validation. In this section, we investigate what patterns have been learned by the model to distinguish between alert and drowsy EEG signals with the visualization technique described in Section III.C. In this connection, we display some representative samples that are correctly classified with high likelihoods of the alert label and the drowsy label in Figure 4 and Figure 5, respectively.

By observation on the relative power of the samples from Figure 4, we find that the selected EEG samples with high likelihood of the drowsy label commonly contain a high portion of Theta waves, e.g., Figure 4(a), or Alpha waves, e.g., Figure 4 (b-d). For the first sample shown in Figure 4(a), it can be observed that the model has identified several episodes that contain rhythmic bursts in the Theta band as strong evidence of the drowsy state. Actually these bursts in the Theta band, or called "drowsy bursts", have been found to frequently appear in EEG signals during drowsiness [55].

For the samples shown in Figure 4(b-d), it can be observed that the model has identified spindle-like structures in Alpha frequency from several episodes of the signal as indicators of drowsiness. Actually, the captured Alpha spindles, which can be characterized by a narrow frequency peak within the alpha band [56], have been found to be strong indicators of early drowsiness in various driving simulator studies and used to identify the driver drowsiness [3].

Another pattern we have observed is that for the samples classified with a high likelihood of alert labels, the central EEG channels, e.g., CPZ in Figure 4(c), usually play a more importance role than peripheral channels for the classification. We infer the reason is that these channels mostly contain cleaner cortical signals where the drowsiness-related features are more distinguishable than that from the peripheral and frontal channels, which are more likely to be contaminated by artifacts caused by brain muscle tension or eye movements.



(a)



(b)

(c)



(d)

**Figure 5**. Visualization of learned patterns on selected alert EEG samples that are correctly classified by the network with high likelihood.

As it can be seen in Figure 5, we have found that the samples classified with a high likelihood of the alert label commonly contain a high portion of artifacts in the signals. For the first sample shown in Figure 5(a), we can see that the model has identified several episodes of the signals from channels of F7, T3 and T5 that contain a high portion of Beta waves as evidence for the classification. Similar to the first sample, the detected discriminative areas of the second sample shown in Figure 5(b) also contain a high portion of Beta waves but from channels of O2, T6 and TP8. Actually, the Electromyography (EMG) activities have the greatest contamination on EEG signals at the periphery of the scalp near the active muscles and the spectra of EMG often have peaks in the Beta frequency range that resemble EEG Beta peaks [57]. Therefore, the high-frequency waves identified by the model in samples from Figure 5(a) and 5(b) could be mostly caused by tension of the scalp muscles. For the samples shown in Figure 5(c) and 5(d), it can be observed that the model has identified several episodes that contain large voltage change of signals from frontal EEG channels as evidence of alertness. These large-amplitude and low-frequency waves, resulting a high power in the Delta frequency band, are caused by eye blinks and eye movement activities happening when the subject is in the alert state.

Actually, it is out of our expectation that the model mostly uses features that are commonly regarded as artifacts contained in EEG rather than the cortical signals as indicators of the alert state. In fact, these artifacts usually dominate the wakeful EEG signals [55], while they are not common in drowsy EEG signals, as it can be seen in samples from Figure 4. It makes sense to some extent that the model uses such features to distinguish alert EEG signals from the drowsy EEG signals.
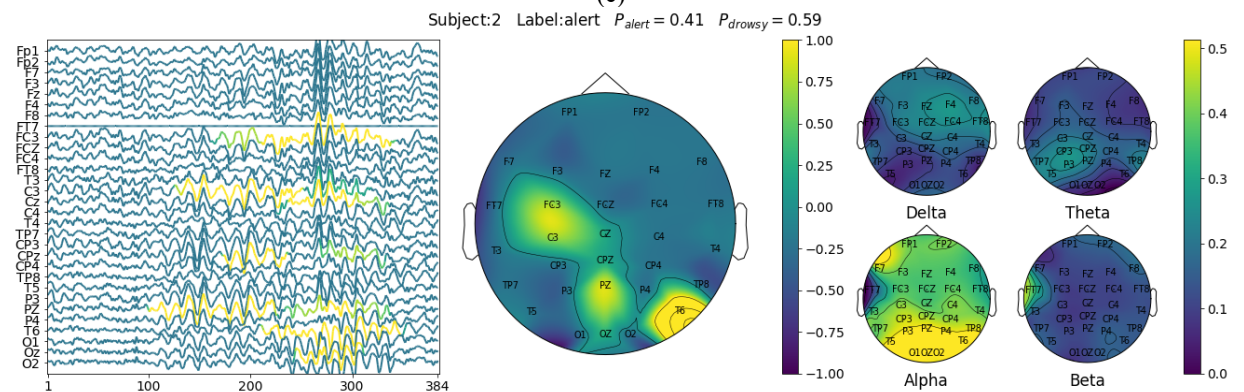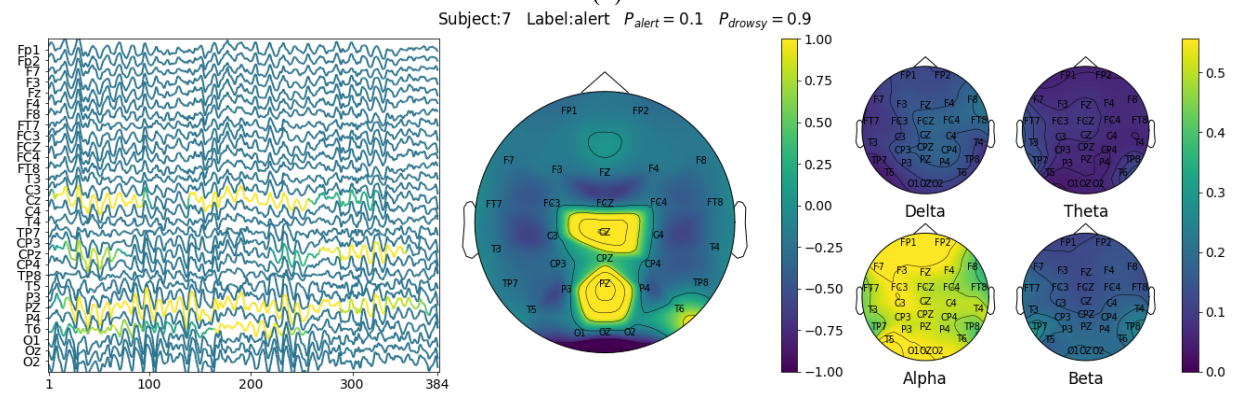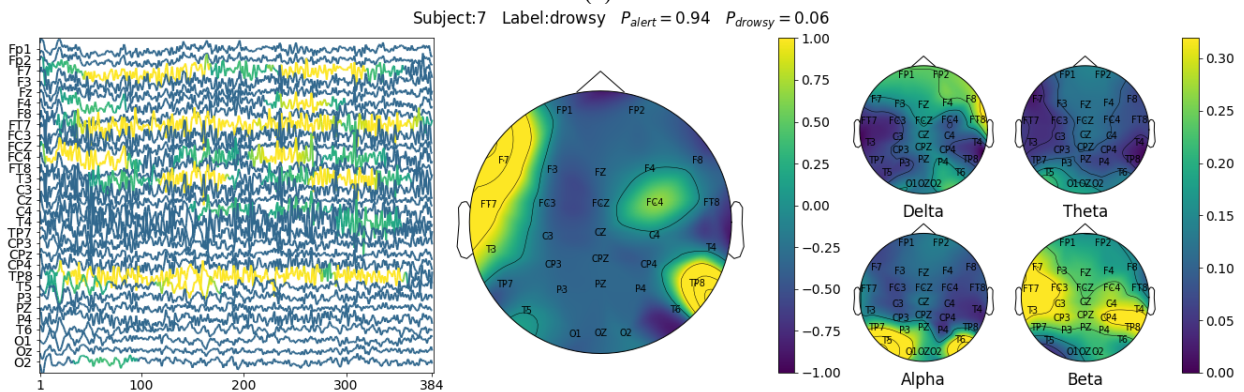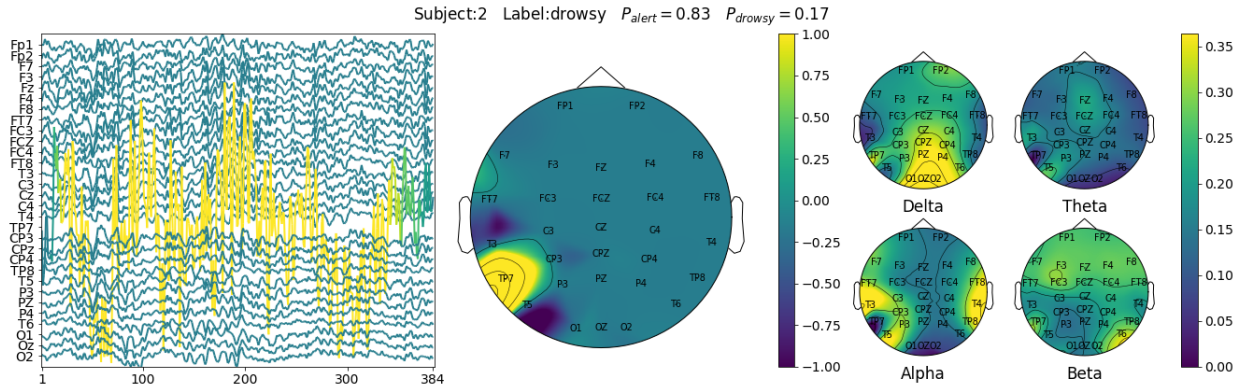
Figure 4. Visualization of learned patterns on selected wrongly classified samples from Subject 2 and Subject 7.

Finally, we consider using the visualization technique to explore the reasons behind some wrongly classified samples. From Table 2, we can see that the proposed CNN model has relatively low classification accuracies of 67.65% and 67.06 for Subject 2 and Subject 7, respectively. Therefore, we visualize some wrongly classified samples from these subjects and display the results in Figure 6.

By analysis of the wrongly classified samples from Subject 2, we find one of the major reasons that lead to the low classification accuracy of the subject is due to the sensor noise contained in the EEG signals. As it can be seen in the example shown in Figure 6(a), the model has falsely identified the noise, which causes significant fluctuations in the channel TP7, as evidence of the alert state. By comparison, the opposite case when the signal is completely lost, as it is shown for the case of channel FT7 from the sample in Figure 6(d), does not affect the classification. The observations above indicate the necessity to filter out the signals with amplitude significantly larger than standard EEG voltage range, in order to prevent them from misleading the model.

The second sample shown in Figure 7(b) contains a high portion of Beta waves, and the model has identified several episodes of signals from the peripheral channels of F7, FT7, T3 and TP8, as evidence of the alert state. Although the sample is labeled with the drowsy state, it does not contain obvious drowsiness-related features, such as alpha spindles or bursts in the Theta-Delta band, for the model to recognize. In fact, the characteristics that are displayed in this sample are quite similar to those of the typical wakeful samples shown in Figure 5(a) and 5(b), and their similarities may explain why the sample is classified with the alert label by the model.

For the last two samples shown in Figure 6(c) and (d), the model has identified several episodes that contain Alpha waves as evidence of drowsiness. For the sample shown in Figure 6(c), it can be seen that the Alpha spindles are generated from the occipital areas with the largest amplitude and propagate to the entire area of the brain. The model has recognized signals from the CZ and PZ channels that contain such spindles as evidence of the drowsy state. In the sample shown in Figure 6(d), the spindles appear from around the $100^{th}$ sampling point and end at around the $350^{th}$ point in almost all the EEG channels. The model has identified several episodes from the central channels and the T6 channel as the evidence for drowsiness. The observations above can justify the classifications by the model, since the Alpha spindles have been found to be strong indicators of drowsiness. In fact, it is possible that the subject was already in the early drowsy stage but coincidentally responded timely when these samples were captured.


## V.  DISCUSSION AND FUTURE WORKS

In this paper, we developed a novel CNN model for the purpose of learning and visualizing common drowsiness-related features in EEG signals across different subjects. The model has a compact structure and it uses separable convolutions to process the EEG signals in a spatial-temporal sequence. In order to allow the model to "explain" its decisions, we designed a visualization technique, which is inspired by the CAM and Fixation-CNN methods, to reveal the important local regions of the sample for the classification. The visualization results indicate that the model has learned biologically explainable features from the EEG signals, e.g., Alpha spindles and Theta bursts, as indicators of the drowsy state. We have also found that different types of artifacts and noise contained in the signals have different impacts on the classification. For example, the EMG and the eye movement artifacts that usually dominate the wakeful EEG signals are identified by the model as evidence of the alert state, while the sensor noise that causes significant fluctuations of the signal negatively affect the classification results. These findings motivate us to use different EEG pre-processing pipelines in our future works to deal with different kinds of artifacts and noise in the signals, according to their impacts on the classifier revealed by the visualization technique.

We have also noticed that for some wrongly classified samples, the model has found valid evidence to justify its classification. Indeed, the performance or behaviors of subjects, e.g., reaction time, may not faithfully reflect mental states of subjects in certain circumstances, which results in bias in the labeling. It could be an interesting topic to incorporate the "network explanation" into the labeling process, instead of merely using thresholds hard-coded on behavior/performance metrics of the subjects.

Currently, we have only tested the proposed model on a public sustained driving dataset with a limited number of samples as an initial attempt. In our future works, we will consider testing the model on more EEG datasets with different volumes. We also consider designing novel visualization techniques that can be applied a boarder range of deep learning models to interpret their classification results.


## VI.  CONCLUSION

In this paper, we developed a novel CNN model for the purpose of discovering common patterns related to different mental states in EEG signals across different subjects. The model has a compact structure and it uses separable convolutions to process the EEG signals in a spatial-temporal sequence. In addition, we also designed visualization techniques to reveal what has been learned by the model for classification by highlighting the relevant parts of the input signal.

Results show that the model achieves an average accuracy of 78.35% for cross-subject drowsiness recognition, which is higher than that of both the conventional baseline methods and the state-of-art deep learning methods. The visualization results show that the model has learned to identify biologically explainable features, e.g., Alpha spindles, from the data and use them as

evidence to distinguish drowsy EEG signals from alert signals. In addition, we also explored reasons behind some wrongly classified samples and investigated how different kinds of noise and artifacts contained in the signals can affect the classification. Our work illustrates a promising direction to use interpretable deep learning models to discover meaningful patterns related to different mental states from complex EEG signals.

## REFERENCES

1.　Torsvall, L., *Sleepiness on the job: continuously measured EEG changes in train drivers.* Electroencephalography and clinical Neurophysiology, 1987. **66**(6): p. 502-511.
2.　Lal, S.K. and A. Craig, *Driver fatigue: electroencephalography and psychological assessment.* Psychophysiology, 2002. **39**(3): p. 313-321.
3.　Simon, M., et al., *EEG alpha spindle measures as indicators of driver fatigue under real traffic conditions.* Clinical Neurophysiology, 2011. **122**(6): p. 1168-1178.
4.　Hu, J. and J. Min, *Automated detection of driver fatigue based on EEG signals using gradient boosting decision tree model.* Cognitive neurodynamics, 2018. **12**(4): p. 431-440.
5.　Luo, H., et al., *Research on fatigue driving detection using forehead EEG based on adaptive multi-scale entropy.* Biomedical Signal Processing and Control, 2019. **51**: p. 50-58.
6.　Krizhevsky, A., I. Sutskever, and G.E. Hinton. *Imagenet classification with deep convolutional neural networks*. in *Advances in neural information processing systems*. 2012.
7.　Simonyan, K. and A. Zisserman, *Very deep convolutional networks for large-scale image recognition.* arXiv preprint arXiv:1409.1556, 2014.
8.　Längkvist, M., L. Karlsson, and A. Loutfi, *Sleep stage classification using unsupervised feature learning.* Advances in Artificial Neural Systems, 2012. **2012**.
9.　Lawhern, V.J., et al., *EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces.* Journal of neural engineering, 2018. **15**(5): p. 056013.
10.　Bashivan, P., et al., *Learning representations from EEG with deep recurrent-convolutional neural networks.* arXiv preprint arXiv:1511.06448, 2015.
11.　Åkerstedt, T., G. Kecklund, and A. Knutsson, *Manifest sleepiness and the spectral content of the EEG during shift work.* Sleep, 1991. **14**(3): p. 221-225.
12.　Corsi-Cabrera, M., et al., *Changes in the waking EEG as a consequence of sleep and sleep deprivation.* Sleep, 1992. **15**(6): p. 550-555.
13.　Klimesch, W., *EEG alpha and theta oscillations reflect cognitive and memory performance: a review and analysis.* Brain research reviews, 1999. **29**(2-3): p. 169-195.
14.　Yeo, M.V., et al., *Can SVM be used for automatic EEG detection of drowsiness during car driving?* Safety Science, 2009. **47**(1): p. 115-124.
15.　Hu, J., *Automated detection of driver fatigue based on AdaBoost classifier with EEG signals.* Frontiers in computational neuroscience, 2017. **11**: p. 72.
16.　LeCun, Y., Y. Bengio, and G. Hinton, *Deep learning.* nature, 2015. **521**(7553): p. 436-444.
17.　Hinton, G., et al., *Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups.* IEEE Signal processing magazine, 2012. **29**(6): p. 82-97.
18.　Schirrmeister, R.T., et al., *Deep learning with convolutional neural networks for brain mapping and decoding of movement-related information from the human EEG. arXiv, 2017.* arXiv preprint arXiv:1703.05051.
19.　Ang, K.K., et al. *Filter bank common spatial pattern (FBCSP) in brain-computer interface*. in *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*. 2008. IEEE.
20.　Roy, Y., et al., *Deep learning-based electroencephalography analysis: a systematic review.* Journal of neural engineering, 2019. **16**(5): p. 051001.
21.　Rundo, F., et al., *An innovative deep learning algorithm for drowsiness detection from EEG signal.* Computation, 2019. **7**(1): p. 13.
22.　Conoci, S., et al. *Advanced skin lesion discrimination pipeline for early melanoma cancer diagnosis towards PoC devices*. in *2017 European Conference on Circuit Theory and Design (ECCTD)*. 2017. IEEE.

23. Nissimagoudar, P.C., A.V. Nandi, and H. Gireesha. *Deep Convolution Neural Network-Based Feature Learning Model for EEG Based Driver Alert/Drowsy State Detection*. in *International Conference on Soft Computing and Pattern Recognition*. 2019. Springer.

24. Ding, S., et al. *Cascaded Convolutional Neural Network with Attention Mechanism for Mobile EEG-based Driver Drowsiness Detection System*. in *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. 2019. IEEE.

25. Li, G., B.-L. Lee, and W.-Y. Chung, *Smartwatch-based wearable EEG system for driver drowsiness detection.* IEEE Sensors Journal, 2015. **15**(12): p. 7169-7180.

26. Gao, Z., et al., *EEG-based spatio–temporal convolutional neural network for driver fatigue evaluation.* IEEE transactions on neural networks and learning systems, 2019. **30**(9): p. 2755-2763.

27. Zeng, H., et al., *EEG classification of driver mental states by deep learning.* Cognitive neurodynamics, 2018. **12**(6): p. 597-606.

28. Hajinoroozi, M., et al., *EEG-based prediction of driver's cognitive performance by deep convolutional neural network.* Signal Processing: Image Communication, 2016. **47**: p. 549-555.

29. Cao, Z., et al., *Multi-channel EEG recordings during a sustained-attention driving task.* Scientific data, 2019. **6**(1): p. 1-8.

30. Lin, Z.C.M.C.J.T.K.C.-T. *Multi-channel EEG recordings during a sustained-attention driving task (pre-processed dataset)*. 2019; Available from: https://figshare.com/articles/Multi-channel_EEG_recordings_during_a_sustained-attention_driving_task_preprocessed_dataset_/7666055.

31. Pal, N.R., et al., *EEG-based subject-and session-independent drowsiness detection: an unsupervised approach.* EURASIP Journal on Advances in Signal Processing, 2008. **2008**(1): p. 519480.

32. Baig, M.Z., N. Aslam, and H.P. Shum, *Filtering techniques for channel selection in motor imagery EEG applications: a survey.* Artificial intelligence review, 2020. **53**(2): p. 1207-1232.

33. Makeig, S., et al. *Independent component analysis of electroencephalographic data*. in *Advances in neural information processing systems*. 1996.

34. Wang, Y., S. Gao, and X. Gao. *Common spatial pattern method for channel selelction in motor imagery based brain-computer interface*. in *2005 IEEE engineering in medicine and biology 27th annual conference*. 2006. IEEE.

35. Friman, O., I. Volosyak, and A. Graser, *Multiple channel detection of steady-state visual evoked potentials for brain-computer interfaces.* IEEE transactions on biomedical engineering, 2007. **54**(4): p. 742-750.

36. Lin, Z., et al., *Frequency recognition based on canonical correlation analysis for SSVEP-based BCIs.* IEEE transactions on biomedical engineering, 2006. **53**(12): p. 2610-2614.

37. Lin, M., Q. Chen, and S. Yan, *Network in network.* arXiv preprint arXiv:1312.4400, 2013.

38. Chollet, F. *Xception: Deep learning with depthwise separable convolutions*. in *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.

39. Lin, C.-T., et al., *EEG-based drowsiness estimation for safety driving using independent component analysis.* IEEE Transactions on Circuits and Systems I: Regular Papers, 2005. **52**(12): p. 2726-2738.

40. Fahimi, F., et al., *Inter-subject transfer learning with an end-to-end deep convolutional neural network for EEG-based BCI.* Journal of neural engineering, 2019. **16**(2): p. 026007.

41. Zhou, B., et al. *Learning deep features for discriminative localization*. in *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

42. Mopuri, K.R., U. Garg, and R.V. Babu, *Cnn fixations: an unraveling approach to visualize the discriminative image regions.* IEEE Transactions on Image Processing, 2018. **28**(5): p. 2116-2125.

43. Ang, K.K., et al., *Filter bank common spatial pattern algorithm on BCI competition IV datasets 2a and 2b.* Frontiers in neuroscience, 2012. **6**: p. 39.

44. Liu, Y., et al. *EEG-Based Cross-Subject Mental Fatigue Recognition*. in *2019 International Conference on Cyberworlds (CW)*. 2019. IEEE.

45. Onton, J., et al., *Imaging human EEG dynamics using independent component analysis.* Neuroscience & biobehavioral reviews, 2006. **30**(6): p. 808-822.

46. Jap, B.T., et al., *Using EEG spectral components to assess algorithms for detecting fatigue.* Expert Systems with Applications, 2009. **36**(2): p. 2352-2359.

47. Wang, Q., Y. Li, and X. Liu, *Analysis of feature fatigue EEG signals based on wavelet entropy.* International Journal of Pattern Recognition and Artificial Intelligence, 2018. **32**(08): p. 1854023.

48. Song, Y., J. Crowcroft, and J. Zhang, *Automatic epileptic seizure detection in EEGs based on optimized sample entropy and extreme learning machine.* Journal of neuroscience methods, 2012. **210**(2): p. 132-146.

49. Xiang, J., et al., *The detection of epileptic seizure signals based on fuzzy entropy.* Journal of neuroscience methods, 2015. **243**: p. 18-25.

50. Farahat, A., et al., *Convolutional neural networks for decoding of covert attention focus and saliency maps for EEG feature visualization.* Journal of neural engineering, 2019. **16**(6): p. 066010.
51. vlawhern, *Army Research Laboratory (ARL) EEGModels Project: A Collection of Convolutional Neural Network (CNN) models for EEG signal classification, using Keras and Tensorflow.*
52. *SciPy.org.* Available from: https://www.scipy.org/.
53. *scikit-lean Machine Learning in Python.* Available from: https://scikit-learn.org/stable/.
54. Kingma, D.P. and J. Ba, *Adam: A method for stochastic optimization.* arXiv preprint arXiv:1412.6980, 2014.
55. Britton, J.W., et al., *Electroencephalography (EEG): An introductory text and atlas of normal and abnormal findings in adults, children, and infants.* 2016.
56. Shaw, J.C., *The brain's alpha rhythms and the mind.* BV Elsevier Science, 2003.
57. Goncharova, I.I., et al., *EMG contamination of EEG: spectral and topographical characteristics.* Clinical neurophysiology, 2003. **114**(9): p. 1580-1593.