# Multi-Scale Neural Network for EEG Representation Learning in BCI

Wonjun Ko, Eunjin Jeon, Seungwoo Jeong, and Heung-Il Suk, *Member, IEEE*

*Abstract*—Recent advances in deep learning have had a methodological and practical impact on brain–computer interface (BCI) research. Among the various deep network architectures, convolutional neural networks (CNNs) have been well suited for *spatio-spectral-temporal* electroencephalogram (EEG) signal representation learning. Most of the existing CNN-based methods described in the literature extract features at a sequential level of abstraction with repetitive nonlinear operations and involve densely connected layers for classification. However, studies in neurophysiology have revealed that EEG signals carry information in different ranges of frequency components. To better reflect these *multi-frequency* properties in EEGs, we propose a novel deep *multi-scale neural network* that discovers feature representations in multiple frequency/time ranges and extracts relationships among electrodes, *i.e.*, spatial representations, for subject intention/condition identification. Furthermore, by completely representing EEG signals with spatio-spectral-temporal information, the proposed method can be utilized for diverse paradigms in both active and passive BCIs, contrary to existing methods that are primarily focused on single-paradigm BCIs. To demonstrate the validity of our proposed method, we conducted experiments on various paradigms of active/passive BCI datasets. Our experimental results demonstrated that the proposed method achieved performance improvements when judged against comparable state-of-the-art methods. Additionally, we analyzed the proposed method using different techniques, such as PSD curves and relevance score inspection to validate the multi-scale EEG signal information capturing ability, activation pattern maps for investigating the learned spatial filters, and t-SNE plotting for visualizing represented features. Finally, we also demonstrated our method's application to real-world problems.

*Index Terms*—Active/Passive Brain–Computer Interface; Electroencephalogram; Deep Learning; Convolutional Neural Network; Motor Imagery; Steady-State Visually Evoked Potentials; Mental Fatigue; Seizure

## I. INTRODUCTION

**B**RAIN–computer interface (BCI) [1] is an emerging technology that enables a communication pathway between a user and an external device (*e.g.*, a computer) through the acquisition and analysis of brain signals. Then these signals are translated into commands that are understood by a device, such as a computer. Owing to its practicality, electroencephalogram (EEG)-based non-invasive BCIs are widely used [1]–[3]. Earlier, Aricò *et al.* [4] categorized user-centered BCIs into two types, active/reactive and passive BCIs. In this paper, our focus is not only on active BCIs but also on passive BCIs. Generally, two types of brain signals such as *evoked* and *spontaneous* EEG are primarily considered for active/reactive

W. Ko, E. Jeon, and H.-I. Suk are with the Department of Brain and Cognitive Engineering, Korea University, Seoul 02841, Korea. S. Jeong and H.-I. Suk are with the Department of Artificial Intelligence, Korea University, Seoul 02841, Korea. Correspondence: hisuk@korea.ac.kr (Heung-Il Suk)

BCIs [5]. Evoked BCIs exploit unintentional electrical potentials reacting to external or internal stimuli. Examples of evoked BCIs include steady-state visually evoked potentials (SSVEP) [6], [7] and event-related potentials [6]. Additionally, spontaneous BCIs use an internal cognitive process such as event related desynchronization and event related synchronization (ERD/ERS) in sensorimotor rhythms, *e.g.*, motor imagery (MI) [6], [8] induced by imagining movements in addition to physical movement. Well-known examples of passive BCIs include the use of sleep/drowsy EEG signals for sleep stage classification or identifying mental fatigue to alert a driver of a dangerous situation and seizure EEG patterns for onset detection to provide the patient with a warning of a potential seizure.

Generally, machine learning-based BCIs consist of five main processing stages [3]: (i) an EEG signal acquisition phase based on each paradigm, (ii) signal preprocessing (*e.g.*, channel selection and band-pass filtering), (iii) feature representation learning, (iv) classifier learning, and finally (v) a feedback stage. Basically, most of machine learning-based BCI methods follow these processes, however, these methods need specific modification to classify a user's intention/condition for each different paradigm [3]. In other words, machine learning-based methods need to have *prior* knowledge of different EEG paradigms [1], [3], [6], [7], [11]. Therefore, conventional machine learning-based BCIs have discovered EEG representations through extremely specialized approaches, *e.g.*, a common spatial pattern (CSP) [1] or its variants [12], [13] for MI signals and a canonical correlation analysis (CCA) [7] for SSVEP signals decoding.

While hand-crafted feature representation learning has a pivotal role in a conventional machine learning framework [1], [7], [14], deep learning-based representation has had remarkable results in the BCI community [2], [3], [5]. These deep learning-based methods have integrated a feature extraction step with a classifier learning step such that those steps are jointly optimized, thereby improving performance. Among various deep learning methods, convolutional neural networks (CNNs) have the advantage [3], [15], [16] of maintaining the structural and configurational information in the original data. In this respect, developing a novel CNN architecture for EEG signal representation has taken a center stage in the BCI studies [2], [15], [17]–[23].

However, some challenges still remain. First, existing CNN-based methods [2], [15], [17], [18], [22], [23] are mostly comprised of stacked convolutional layers. In other words, those existing methods extract features sequentially. But, ignoring multiple ranges of spectral-temporal features can cause
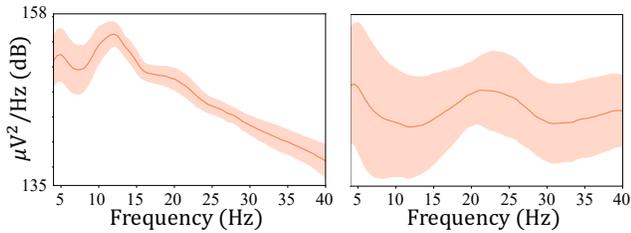
Fig. 1: Power spectral density (PSD) curves of two different subjects' MI EEG samples. The solid red line denotes the mean PSD and the shaded region exhibits the standard deviation for all trials. Clearly, these two different subjects show quite different PSD patterns for the same paradigm (motor imagery).

a critical problem because EEG signal features for different subjects [24], paradigms [3], and types [4] are found in diverse ranges. For example, Fig. 1 depicts two different subjects' MI EEG power spectral density (PSD) curves. Clearly, these two plots have different distributions from each other even though these PSDs are estimated by the same task. Therefore, it is important to capture multi-scale spectral information in EEGs for *general use in* BCI, *i.e.*, a generic method applicable to various types of BCIs.

In addition, those stacked CNN-based methods [2], [15], [22], [23] have numerous trainable parameters, thus requiring large amounts of training samples, whereas BCIs generally acquire a limited number of EEG trials [24]. Therefore, generalizing conventional stacked CNN-based methods in BCI is quite difficult because deep learning is a *data-hungry* problem, *i.e.*, rarely generalized with a lack of data.

Finally, interpreting a learned stacked CNN from a neurophysiologically appropriate standpoint [25] is quite complicated because the CNN identifies complex patterns of data in latent space making a direct explanation difficult [25].

In this study, we propose a novel deep learning-based BCI method to mitigate the previously discussed difficulties. The main contributions of our study are as follows:

- First, we propose a novel CNN architecture, that is applicable independently from the input paradigm or type of EEG and can represent multi-scale spatio-spectral-temporal features.
- Second, the proposed method achieved positive performance on five different datasets for four differnt paradigms (two for active BCIs and two for passive BCIs). The proposed method outperformed or was similar to state-of-the-art linear and deep learning methods, which were individually designed for each specific paradigm.
- Last, we analyze the proposed network using a variety of techniques.

The rest of this paper is organized as follows: Section II reviews previous research on various EEG representation learning via linear model-based or deep learning-based methods. In Section III, we propose a novel and compact deep CNN that classifies multi-paradigm EEG by representing multi-scale spatio-spectral-temporal features. Section IV presents experimental settings and results by comparing the proposed method and comparable baselines. In Section V, we analyze our proposed method from several points of view. Finally, Section VI summarizes the proposed study and suggest future research directions.

## II. RELATED WORK

Learning a class-discriminative feature representation of EEG is still challenging in both theory and practice. Numerous prior studies have attempted to extract features from EEGs. In this section, we briefly discuss linear methods and deep learning models used for EEG signal representation.

### A. Linear Models

Over the past decades, CSP [1] and its variants [12], [13] have played an essential role in decoding MI. Blankertz *et al.* [1] and Ang *et al.* [13] independently used a spatial filtering-based method for classifying MI. Ang *et al.* [13] band-pass filtered EEG data before applying CSP, thereby attempting to decode EEG signals in a spatio-spectral manner. They named the proposed method filter bank CSP (FBCSP). Furthermore, Suk and Lee [12] also decoded MI by jointly optimizing multi spectral filters in a Bayesian framework.

CCA is commonly utilized for detecting SSVEP [7] owing to its practical ability to be implemented without the calibration stage. The standard CCA method [7] deployed sinusoidal signals as reference signals and estimated canonical correlation between the reference signals and input EEG signals to identify an evoked frequency in SSVEP EEGs.

In addition, to characterize the sleep stage, entropy calculation-based approaches were frequently used. Sanders *et al.* [26] classified the sleep stage using the spectral-temporal features of EEGs learned from short-time Fourier transformation. Furthermore, Zheng and Lu [10] focused on identifying a driver's mental fatigue during driving. They [10] applied filter banks to EEG signals to extract spectral information, and then transformed the filtered EEG signals to spectral space, *i.e.*, estimated PSD of filtered EEG signals. By doing so, Zheng and Lu [10] effectively assessed the regression score of the driver's mental states which were labeled using the PERCLOS index, a measure of neurophysiological fatigue.

Earlier, Shoeb and Guttag [14] applied a machine learning approach to extract and classify the spatio-spectral-temporal features of epileptic seizure EEG signals. Specifically, these authors [14] used filter banks in a channel-wise manner to capture the spatio-spectral information. Then, by encoding the temporal evolution of extracted spatio-spectral feature vectors, they [14] effectively constructed epileptic seizure EEG signal spatio-spectral-temporal features and classified the seizure and non-seizure features utilizing a support vector machine (SVM). Recently, spectral features derived from a principal component analysis (PCA) [27] exhibited superior performance for seizure onset detection. In particular, Lee *et al.* [27] band-pass filtered raw signals and calculated PSD. Then they [27] applied PCA for the extraction of EEG signal spectral features.

These practical linear model-based BCI methods [1], [7], [10], [12], [13], [26], [27] have demonstrated credible performance. However, these methods need to have certain prior neurophysiology knowledge [3], because their feature extraction stages are specifically designed for each EEG paradigm. Conversely, our method does not need to be specialized for different paradigms.

### B. Deep and Hierarchical Models

Recently, deep learning methods, especially CNNs have achieved promising results in EEG signal decoding researches. For instance, Schirrmeister *et al.* [2] introduced Shallow ConvNet, Deep ConvNet, Hybrid ConvNet, and Residual ConvNet. These authors [2] evaluated how well various proposed CNNs decoded MI. Ko *et al.* [15] also proposed a novel CNN architecture which is inspired by a recurrent convolutional neural network [28] for MI classification, deep recurrent spatio-temporal neural network (RSTNN).

While a standard CCA [7] has obtained state-of-the-art performance in SSVEP BCI, Kwak *et al.* [20] developed a CNN for SSVEP feature representation learning. These authors simply combined spatial and temporal convolution to enable the system to learn data patterns in the latent space, thereby correctly generalizing EEG signal features. Meanwhile, Waytowich *et al.* [19] applied EEGNet [3] to the SSVEP paradigm and achieved a higher performance than that of the standard CCA [7].

Supratak *et al.* [17] developed a deep neural network for sleep stage detection. More precisely, they combined a CNN for representation learning and a recurrent neural network for sequential residual learning [17]. Furthermore, they trained the deep learning model in two separate steps, optimizing the model by individual pre-training and fine-tuning. In the meantime, Gao *et al.* [23] proposed an EEG-based spatio-temporal convolutional neural network (ESTCNN) for driver fatigue evaluation. The ESTCNN [23] simply convolved the band-pass filtered EEG to represent temporal dependencies and flattened the extracted features for spatial features fusion. Lastly, densely connected layers were used for the identification of a user's condition [23].

To detect a seizure type, Asif *et al.* [21] proposed a multispectral deep feature learning using a deep CNN, SeizureNet. These authors [21] transformed the EEG signals to spectral space using saliency-encoded spectrogram generation and fed the extracted spectral features to a deep neural network. In the meantime, Emami *et al.* [22] independently proposed another CNN-based approach for detecting seizure onset. They [22] band-pass filtered and segmented the input EEG patterns, and then used a deep CNN for classification.

Recently, Lawhern *et al.* [3], [19] proposed a novel CNN, *so-called* EEGNet. Unlike other linear or deep learning-based methods, the EEGNet classified various EEG paradigms using a single architecture, *i.e.*, not specifically tuned for different paradigms. Further, Lawhern *et al.* [3] introduced a separable convolution [16] and used it as a parameter reduction method.

On the one hand, the deep and hierarchical models decoded the EEG signals well without any custom feature extraction
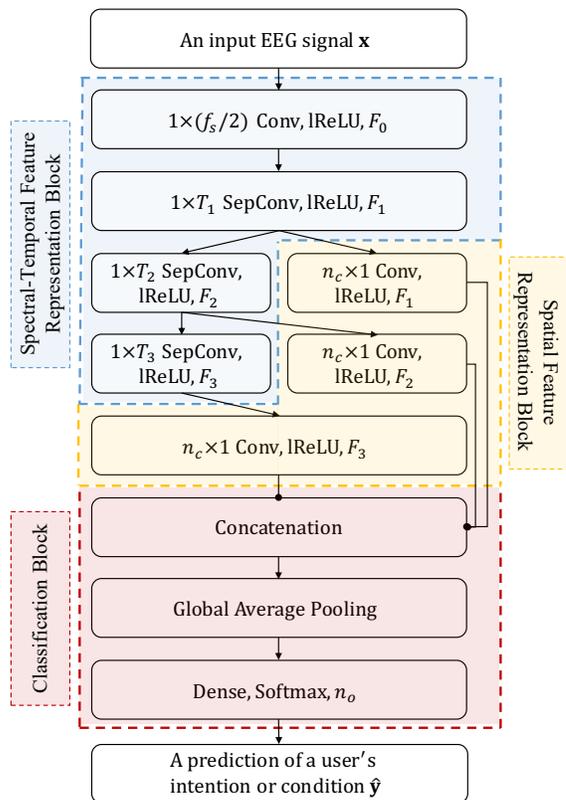


Fig. 2: Architectural framework of our multi-scale neural network (MSNN). In the proposed network, first, an input EEG **x** is temporally convolved to expand the number of features, where $f_s$, $F_0$, lReLU denote the sampling frequency rate, the number of output filter maps of the first layer, and leaky rectified linear unit activation function respectively. Then, a set of temporal separable convolutions extracts *spectral-temporal* features ($T_k$ and $F_k$ respectively denote the kernel size and output feature maps of $k$-th temporal separable convolution). At the same time, a set of spatial convolutions represents *spatial* features, where $n_c$ denotes the number of acquired EEG channels. Then, the *multi-scale* features are concatenated and fed into the global average pooling layer [29]. Finally, the dense layer classifies the class of input EEG by exploiting multi-scale features where $n_o$ denotes the number of output nodes.

stage for their respective paradigm [2], [15], [17], [20]–[23] or even various paradigms [3], [19]. On the other hand, the deep CNNs extracted the EEG features at a sequential level using stacked convolutional layers without exploiting multi-scale spectral representation. Conversely, the proposed method exploits multi-scale spatio-spectral-temporal features irrespective of the input EEG paradigms.

## III. METHODS

In this section, we propose a deep multi-scale neural network (MSNN), which can represent EEG features from different paradigms by exploiting *spatio-spectral-temporal* information at multi-scale.

### A. Multi-Scale Neural Network

As mentioned previously, an FBCSP [13] is one of the most successful models to exploit EEG signal *multi-scale* features, especially for MI. Thus, many successful MI EEG signal decoding algorithms [2], [12] or even other paradigm classification algorithms [3] are inspired by the FBCSP [13] model. In this study, the proposed multi-scale neural network (MSNN) also learns multi-scale feature representations. However, the network automatically learns from data through discriminative multiple spectral filters, rather than manually defining multi-frequency bounds as in FBCSP [13]. Basically, our proposed method consists of three types of blocks: (1) a *spectral-temporal* feature representation block, (2) a *spatial* feature representation block, and, (3) a classification block, as depicted in Fig 2.

First, in the spectral-temporal feature representation block, stacked convolutional layers extract EEG data spectral-temporal features, such as existing EEG classification methods. However, the proposed model exploits intermediate activations for gathering multi-scale spectral information. Then, the spatial feature representation block discovers spatial patterns from the extracted multi-scale features. Finally, these multi-scale spatio-spectral-temporal features are concatenated, pooled, and fed into the densely connected layer for classification.

### B. Spectral-Temporal Feature Representation Block

Given an input EEG data $\mathbf{x}$, we reshape it in the form of $[n_c, n_T, 1]$, *i.e.*, $x \in \mathbb{R}^{n_c \times n_t \times 1}$, where $n_c$ and $n_T$ denote the number of electrode channels and timepoints, respectively.

In the MSNN, the input EEG data are temporally convolved in a channel-wise manner by a temporal convolutional layer to expand the number of feature maps. Thus, the activated features have the form $[n_c, n'_T, F_0]$, where $n'_T = n_T - (f_s/2) + 1$ ($f_s$ and $F_0$ are the sampling frequency and the feature map dimension for the first temporal convolution layer.). The main benefits of using a separable convolution [3], [16] are a significant reduction of tunable weights in the model and, more importantly, an efficient and explicit decoupling of the relationship between the temporal and the feature map dimensions of the input features. This is accomplished by learning kernels independently for each feature map. Thus, as in BCI literature, the separable convolution [16] enables the system to learn temporal kernels individually from the feature map dimensions (using a depthwise convolution [16]), and then optimally re-combine the feature maps (using a pointwise convolution [16]).

In this block, by setting a kernel size of $(1 \times T_k)$, where $T_k$ denotes the kernel size of the $k$-th temporal separable convolution, the $k$-th temporal separable convolutional layer represents EEG signal features in the range of $T_k/f'_s$ sec, hence, $f'_s/T_k$ Hz, where $f'_s$ is a frequency property extracted at the first temporal convolutional layer. Therefore, the spectral-temporal feature representation layers can deal with different timepoints or frequency ranges by using various kernel sizes for the input EEG data.

Additionally, each different layer that has a different kernel size extracts features in different frequency and timepoint ranges. In other words, a spectral-temporal convolution layer with a larger kernel represents *longer-term* temporal features, *i.e.*, a *lower-range* of spectral features and vice versa. Then, the MSNN exploits intermediate activations from each layer, thus learning multi-scale feature representations.

In addition, a separable convolution [16] only operates convolutions in a cross-$n_T$ way, thus, the number of parameters is small compared to a conventional convolution. For instance, while a $k$-th separable temporal convolution has only $T_k + F_{k-1} \cdot F_k$ parameters, the conventional convolution with the same size kernel has $T_k \cdot F_{k-1} \cdot F_k$ parameters, where $F_k$ denotes the feature maps dimension of $k$-th layer.

Furthermore, in this processing, as described above, the MSNN uses its intermediate activations to exploit multi-scale representations. In other words, the proposed network obtains $N$ numbers of spectral-temporal features $\mathbf{f}_k^{\mathrm{ST}}$, $k = 1, 2, \cdots, N$ like:

$$\mathbf{f}_k^{\mathrm{ST}} = \mathcal{C}_k^{\mathrm{ST}}(\mathbf{x}) = \mathcal{C}_k^{\mathrm{sep}} \circ \mathcal{C}_{k-1}^{\mathrm{sep}} \circ \cdots \circ \mathcal{C}_0(\mathbf{x}), \qquad (1)$$

where $\mathcal{C}_k^{\mathrm{sep}}$, $\mathcal{C}_0$, and $\mathcal{F}_i \circ \mathcal{F}_j$, respectively, denote the $k$-th separable convolution, the first temporal convolution, and a function composition between arbitrary functions $\mathcal{F}_i$ and $\mathcal{F}_j$, *i.e.*, $\mathcal{F}_i \circ \mathcal{F}_j(\cdot) = \mathcal{F}_i(\mathcal{F}_j(\cdot))$. Thus, by extracting features $\mathbf{f}_1^{\mathrm{ST}}, \mathbf{f}_2^{\mathrm{ST}}, \cdots, \mathbf{f}_N^{\mathrm{ST}}$, the MSNN effectively represents the spectral-temporal features from the multi-scale viewpoint, thereby automatically enhancing generalization. In addition, as all inputs are *zero-padded* before each separable temporal convolution, the output features have the same dimension for the channels and timepoints, except for the feature map dimension. Thus, the $k$-th spectral-temporal feature $\mathbf{f}_k^{\mathrm{ST}}$ now has the form $[n_c, n'_T, F_k]$.

### C. Spatial Feature Representation Block

In the spatial feature representation block, a common spatial convolution is used for feature extraction. In this block, the kernel size is constrained to be equal to the number of EEG channels, hence, a convolution with a kernel of $(n_c \times 1)$ is used. Additionally, by setting the kernel size to be the same as the number of electrode channels, similar to many existing deep learning-based BCI methods [2], [3], [15], the proposed MSNN extracts spatial information from the original EEG acquisition channel distributions of multi-scale spectral temporal features. Then, the MSNN can obtain *neurophysiologically plausible* information from the input data distribution.

Furthermore, the spatial feature representation can be applied unrestrictedly, thus in the proposed method, we add this block after every extracted spectral-temporal features $\mathbf{f}_k^{\mathrm{ST}}$, $k = 1, 2, \cdots, N$ like,

$$\mathbf{f}_k^{\mathrm{SST}} = \mathcal{S}_k(\mathbf{f}_k^{\mathrm{ST}}) = \mathcal{S}_k \circ \mathcal{C}_k^{\mathrm{ST}}(\mathbf{x}), \qquad (2)$$

where $\mathcal{S}_k$ denotes the $k$-th spatial convolution and $\mathbf{f}_k^{\mathrm{SST}}$ is *spatio-spectral-temporal* features estimated by the $\mathcal{S}_k$ and $\mathcal{C}_k^{\mathrm{ST}}$. We use *valid paddings* for every spatial convolution, thus the $k$-th spatio-spectral-temporal feature $\mathbf{f}_k^{\mathrm{SST}}$ has the form $[1, n'_T, F_k]$. By setting the number of spatial convolutions to be identical to the number of spectral-temporal convolutions,

unlike many previous researches using deep learning for BCI [2], [3], [15], [17], [21], we extract spatial features of each range from spectral-temporal features. In other words, unlike many previous stacked CNNs, the proposed architecture uses every intermediate activated feature set to exploit spatial information, thereby creating the capability to extract various ranges of EEG features at multi-scale.

### D. Classification Block

For classifier learning, because we have $N$ numbers of different (or same when $F_1 = F_2 = \cdots = F_N$) size of spatio-spectral-temporal features $\mathbf{f}_k^{\text{SST}}$, $k = 1, 2, \cdots, N$, the classifier in the proposed method has to concatenate the features in the feature map dimension. Thus, the concatenated feature $\mathbf{f}_{\text{concat}}^{\text{SST}}$ is represented as:

$$\mathbf{f}_{\text{concat}}^{\text{SST}} = \bigg\|_{i=1}^{N} \mathbf{f}_i^{\text{SST}} = \bigg\|_{i=1}^{N} \mathcal{S}_i \circ \mathcal{C}_i^{\text{ST}}(\mathbf{x}), \tag{3}$$

where $\|$ denotes the concatenation operation.

For the classifier network, let us assume that the number of output classification nodes is denoted by $n_o$ and we use a single linear mapping layer. Then, we need to train the large number of $\sum_{i=1}^{N} n_o \cdot n_T' \cdot F_i$ parameters (note that we disregard the bias term for a convenient calculation) because $\mathbf{f}_{\text{concat}}^{\text{SST}}$ has the form $[1, n_T', \sum_{i=1}^{N} F_i]$, and it would still require a large number of training samples. Therefore, after representing the input EEG data multi-scale spatio-spectral-temporal features, the proposed MSNN has one extra operation for reducing the trainable weights. Unlike the existing deep learning-based BCI methods [2], [3], [15], [17], [20], [21], global average pooling (GAP), which is widely used in the computer vision field [29] is performed.

The GAP layer [29], a type of pooling layer, averages nodes from each feature map, thus eliminating the requirement for any window size or stride. By applying GAP [29], our proposed MSNN efficiently extracts significant features. From the BCI literature, the GAP layer [29] can be understood to be a method that can emphasize an important frequency range and its surrounding area for each feature map dimension. Thus, for the extracted multi-scale features in the MSNN, the GAP layer [29] stresses the crucial spectral-temporal part resulting in concise information for the final decision making.

Additionally, the GAP layer [29] significantly reduces the number of classifier parameters used in the proposed MSNN. Specifically, after the GAP layer $\mathcal{G}(\cdot)$, the extracted feature is reduced to the form $[1, 1, \sum_{i=1}^{N} F_i]$, whereas the feature without GAP has the form $[1, n_T', \sum_{i=1}^{N} F_i]$. Therefore, we drastically suppress the trainable parameters in the classifier from $n_T' \cdot n_o \cdot \sum_{i=1}^{N} F_i$ to $n_o \cdot \sum_{i=1}^{N} F_i$.

Then, the MSNN prediction, $\hat{\mathbf{y}}$, for the input EEG data, $\mathbf{x}$, is as follows:

$$\hat{\mathbf{y}} = \text{softmax}\left(\mathbf{W}_o^{\top} \cdot \mathcal{G}\left(\mathbf{f}_{\text{concat}}^{\text{SST}}\right) + \mathbf{b}_o\right)$$
$$= \text{softmax}\left(\mathbf{W}_o^{\top} \cdot \mathcal{G}\left[\bigg\|_{i=1}^{N} \mathcal{S}_i \circ \mathcal{C}_i^{\text{ST}}(\mathbf{x})\right] + \mathbf{b}_o\right), \tag{4}$$

where $\mathbf{W}_o \in \mathbb{R}^{\sum_{i=1}^{N} F_i \times n_0}$ and $\mathbf{b}_0 \in \mathbb{R}^{n_o}$ respectively denote the weight matrix and bias of the classifier.

Finally, the cross-entropy loss, $\mathcal{L}$, that is used for network training is calculated by the prediction $\hat{\mathbf{y}}$ and the label $\mathbf{y}$:

$$\mathcal{L} = \text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{b=1}^{B} \mathbf{y}^{(b)} \log \hat{\mathbf{y}}^{(b)}, \tag{5}$$

where $B$ and $\text{CE}(\cdot, \cdot)$ respectively denote the mini-batch sizes and the cross-entropy loss function, and $\hat{\mathbf{y}}^{(b)}$ and $\mathbf{y}^{(b)}$ denote the prediction and ground-truth label for the $b$-th training sample in the mini-batch[1].

## IV. EXPERIMENTS

In this section, we describe the datasets used for performance evaluation, our experimental settings, and baseline settings. Furthermore, we present the performance of our method and competing methods.

### A. Datasets and Preprocessing

In this study, we used five different publicly available datasets to validate the proposed method on four different EEG data paradigms.

*1) Motor Imagery:* First, we used two big datasets for MI EEGs, GIST-MI [8][2] and KU-MI [6][3,4]. The GIST-MI [8] dataset consists of two different MI tasks: left-hand and right-hand MI that are acquired from 52 subjects. All EEG signals were recorded from 64 Ag/AgCl electrode channels according to the standard 10-20 system, sampled at 512Hz. Each class contained 100 or 120 trials, and each trial was a 3 sec long MI task. Because this dataset is not separated into training and test samples, we conducted a five-fold cross-validation for a fair evaluation. For the MI datasets, we preprocessed signals by applying a large Laplacian filtering[5], baseline correction by subtracting the mean value of the fixation signal from each MI trial, and band-pass filtering between 4 and 40Hz. Then, we removed the first and last 0.5 sec from each trial, and finally applied Gaussian normalization. We applied the same mean and standard deviation values for normalization to the test samples. The multi-channel EEG signals were only shifted and scaled by their respective channel-wise mean and standard deviation values. Thus, inter-channel relations inherent in the data were preserved.

*2) Steady-State Visually Evoked Potentials:* We also used the KU-SSVEP dataset [6][3] for SSVEP decoding experiments in this study. This KU-SSVEP dataset [6] was acquired from 54 subjects and recorded from 62 Ag/AgCl electrode channels using the 10-20 system. The KU-SSVEP dataset [6] contains four EEG classes from target stimuli at 5.45, 6.67, 8.57,

---

[1] All codes used in our experiments are available at 'https://github.com/DeepBCI/Deep-BCI/tree/master/1_Intelligent_BCI/Multi_Scale_Neural_Network_for_EEG_Representation_Learning_in_BCI.'

[2] Available at http://gigadb.org/dataset/100295

[3] Available at http://gigadb.org/dataset/100542

[4] Experimental results of the KU-MI dataset [6] are reported in Supplementary B.

[5] When the target channel does not have four nearest neighbors, we just used available channels and their average value to filter the target channel.

and 12Hz, and each class has 25 EEG trials of training and testing samples for each session. We preprocessed the SSVEP signals by applying band-pass filtering between 4 and 15Hz and selected eight channels in the occipital region, 'PO3, POz, PO4, PO9, O1, Oz, O2, and PO10,' because this region is widely used for SSVEP classification [19].

*3) Drowsiness:* With respect to passive BCI [4], we considered two different paradigms, seizure EEG signals [11] and vigilance EEG signals [10]. Owing to its theoretical and practical benefits, in this study, we conducted experiments identifying drivers' mental fatigue. We also used a publicly available SEED-VIG EEG dataset [10][6] for the drowsy driving task data. This dataset [10] consists of 23 experiments, *i.e.*, trials, and each trial is recorded for approximately 2 hours while simulated driving occurs. The EEG signals are acquired from 17 electrode channels according to the 10-20 system and sampled at 200Hz [10]. For this dataset, we band-pass filtered EEG signals in the range between 0.5 and 40Hz, each epoch was 8 sec in length. Because the dataset was originally labeled using *PERCLOS* levels [10], we categorized the label vectors into three classes, *awake*, *tired*, and *drowsy* with two threshold values(0.35 and 0.7) [10]. Then, for every 23 experiments, a five-fold cross-validation was used for performance estimations.

*4) Seizure:* Finally, we conducted seizure onset detection experiments with the widely used and publicly available CHB-MIT [11][7] dataset. The CHB-MIT dataset [11] contains EEG data from 24 subjects sampled at 256Hz acquired from 23 electrode channels (24 or 26 in a few cases) according to the 10-20 system. In this work, we selected EEG trials that have the same 23 channels montage and removed some trials acquired from the different montage. By following [14], we used a *leave-one-record-out* cross-validation. More precisely, we trained the proposed method using all non-seizure records and all seizure records but one, and tested the model on the remaining seizure record [14]. Then, we repeated this process for the number of seizure records in the dataset, thus, each seizure record was tested. For training, the test trial epochs were 10 sec in length. During validation and testing session, a 10 sec length EEG signal was input into the proposed network using a 1/256 stride. Then, we observed whether the probability values for each EEG signal timepoint was ictal or normal.

For all datasets, the training samples were randomly selected and split again into training and validation samples for model selection. Specifically, we divided the training samples at a 9:1 ratio for each subject and used them for training and model selection respectively.

## B. Experimental Settings

In our work, we compared our method with paradigm-specific linear model-based and deep learning-based methods for each EEG paradigm.

*1) Linear Models - Motor Imagery:* First, we built a CSP with a linear discriminant analysis (CSP + LDA) [1] and an FBCSP with an LDA (FBCSP + LDA) [13] for MI decoding. We used four filters and regularized covariance for the CSP [1] and FBCSP [13]. Additionally, we also used nine non-overlapped filter banks in the 4∼40Hz range, *i.e.*, 4∼8, 8∼12, ⋯, 36∼40Hz, and, finally selected 10 features using the mutual information-based feature selection method FBCSP [13].

*2) Linear Models - Steady-State Visually Evoked Potentials:* We also built a standard CCA [7] for SSVEP classification. We set reference signals for each stimulus including second harmonics. Furthermore, the standard CCA [7] does not require training samples for the optimization, thus we only estimated each session in its entirety from the KU-SSVEP dataset [6] for the CCA performance estimation.

*3) Linear Models - Drowsiness:* For the drowsy state detection experiment, we estimated the filter-banked input EEG data PSD in a channel-wise manner for extracting spatio-spectral features and classified the learned features using an SVM with a radial basis function (RBF) kernel ($\gamma = 1/d_{\text{input}}$ where $d_{\text{input}}$ denotes the input feature dimension) [10].

*4) Linear Models - Seizure:* In addition, we also reimplemented Shoeb and Guttag [14]'s method for the seizure onset detection experiment. We applied the PSD to the EEG data in a channel-wise manner. Then, the 3 sec time window time evolution [14] method was used for capturing temporal information. Finally, the represented spatio-spectral-temporal features were fed into an SVM using an RBF kernel ($\gamma = 1/d_{\text{input}}$).

*5) Deep Neural Networks - Motor Imagery:* We also implemented deep learning-based BCI models[8] for MI. Basically, most of the existing deep learning models [2], [15], [22], [23] have focused on a paradigm-specific BCI task. However, we conducted experiments over all types of datasets for each deep learning model to demonstrate the validity of the proposed method. We built a Shallow ConvNet and a Deep ConvNet as proposed by Schirrmeister *et al.* [2]. The Shallow ConvNet [2] consists of two convolutions, temporal and spatial, with a squaring nonlinear activation, an average pooling, and a logarithmic activation. The Deep ConvNet [2] has five convolutions, temporal and spatial, and three additional temporal convolutions. The RSTNN [15] is also used for these experiments. This network [15] consists of three recurrent convolutional layers, and each recurrent convolutional layer has three recurrent temporal convolutions [28] with a spatial convolution.

*6) Deep Neural Networks - Steady-State Visually Evoked Potentials:* For the SSVEP decoding experiment, we exploited another version of EEGNet for SSVEP EEG [19]. We used different kernel sizes for this EEGNet [19] as Waytowich *et al.* proposed. The SSVEP classification performance estimated by this version [19] is marked by † in the classification table.

*7) Deep Neural Networks - Drowsiness:* The ESTCNN [23] which is proposed for mental fatigue classification has three core blocks. Each block in the ESTCNN [23] consists of

---

TABLE I: Performance evaluations. The Method column denotes all used classification/detection methods including baselines and the proposed method on the various datasets, GIST-MI [8], KU-SSVEP [6], SEED-VIG [10], and CHB-MIT [11] EEG dataset. Each cell depicts the average performance and the standard deviation of all subjects (or trials for the SEED-VIG [10]). For classification performance on the SSVEP dataset, we used different kernel sizes for EEGNet [19] and the proposed method. These values are marked by † and ‡, respectively.

| Method | GIST-MI [8] | KU-SSVEP [6] | SEED-VIG [10] | | CHB-MIT [11] |
| --- | --- | --- | --- | --- | --- |
| | Classification accuracy | | Number of false detections | | |
| | Mean±Std. | | Mean±Std. | False Positive (Drowsy) | Mean (Mean latency) |
| CSP + LDA [1] | .66±.14 | - | - | - | - |
| FBCSP + LDA [13] | .68±.15 | - | - | - | - |
| CCA [7] | - | .94±.09 | - | - | - |
| PSD + SVM [10] | - | - | 31.20±15.47 | 6.74 | - |
| Shoeb and Guttag [14] | - | - | - | - | 5.35 (5.11) |
| Shallow ConvNet [2] | .63±.11 | .52±.20 | 34.89±19.13 | 6.51 | 19.21 (8.48) |
| Deep ConvNet [2] | .61±.07 | **.96±.08** | 41.31±21.04 | 8.65 | 8.74 (7.52) |
| RSTNN [15] | .69±.12 | .65±.20 | 39.84±22.56 | 8.08 | 24.35 (9.31) |
| ESTCNN [23] | .67±.10 | .79±.17 | 41.10±21.31 | 8.71 | 6.41 (7.01) |
| EEGNet [3], [19] | .64±.07 | .93±.10$^{\dagger}$ | 46.63±22.10 | 11.26 | 5.40 (6.23) |
| MSNN (Proposed) | **.81±.12** | .93±.08$^{\ddagger}$ | **31.10±17.29** | **5.38** | **5.35 (4.98)** |

three temporal convolutions with a max pooling layer with the exception of the last block that uses an average pooling layer instead of the max pooling.

*8) Deep Neural Networks - Multi-paradigm:* Finally, we also implemented the EEGNet [3] in our study. As previously mentioned, we used different kernel sizes for two different EEGNets, [3] and [19]. Nevertheless, the basic architecture of the network was the same for various EEG paradigms, having a temporal convolution, depthwise spatial convolution [16], and separable temporal convolution [16].

*9) Proposed Multi-Scale Neural Network:* While training our proposed network, depicted in Fig. 2, we set a mini-batch size of 16, an exponentially decreasing learning rate (initial value: 0.03, decreasing ratio per epoch: 0.001), and an Adam optimizer. For the first temporal convolution, we used a conventional temporal convolution with the kernel size of $(1 \times f_s/2)$ and $F_0 = 4$. Furthermore, we used three spectral-temporal feature representation convolutions, *i.e.*, $N = 3$, and set $T_1 = 100$, $T_2 = 60$, and $T_3 = 20$ with $F_1 = 16$, $F_2 = 32$, and $F_3 = 64$. Then, for the spatial feature representation block, we used three spatial convolutions because the number of spatial convolutional layers must be the same as the number of spectral-temporal separable convolutional layers. The proposed method used different kernel sizes for the SSVEP dataset, similar to the EEGNet [19] due to the fact that SSVEP EEG data is created by target frequencies [6], [7]. For the KU-SSVEP dataset [6], we set $T_1 = 20$, $T_2 = 10$, and $T_3 = 5$ for the spectral-temporal feature representation block, and used the same settings for the others. The SSVEP classification performance estimated by this method is marked by ‡. Additionally, batch normalization was performed after every convolution. Finally, for the classification block, all activated features from the *spatio-spectral-temporal* block were concatenated and fed into the GAP [29] layer. Then, after flattening, the *multi-scale* features were linearly mapped by a dense layer. In this proposed network, a leaky rectified linear unit (ReLU) activation function, an L1-L2 regularizer ($\ell_1 = 0.01$ and $\ell_2 = 0.001$), and a Xavier initializer [30] are used for all tunable parameters except for the final

decision layer that is activated by a softmax activation function instead of a leaky ReLU. We selected model components that demonstrated the best performance for validation, *i.e.*, model selection samples, as mentioned previously.

### C. Experimental Results

*1) Motor Imagery:* All experimental results are summarized in TABLE I. Our proposed network clearly outperformed other baselines for MI EEG signal decoding. Importantly, the proposed network achieved a higher accuracy than those methods designed specifically for MI classification: CSP [1], FBCSP [13], Shallow ConvNet [2], Deep ConvNet [2], and RSTNN [15]. With this clear improvement in accuracy, we could expect that our proposed method is one step closer to MI-based BCI commercialization.

*2) Steady-State Visually Evoked Potentials:* Our proposed MSNN achieved a slightly lower performance than CCA [7], Deep ConvNet [2], and EEGNet [19] in the SSVEP classification. However, the difference in performance between our MSNN and the other three baselines, CCA [7], Deep ConvNet [2], and EEGNet [19], was reasonably small and the proposed method performed with a credible accuracy score.

*3) Drowsiness:* The proposed MSNN made the smallest number of mistakes in decision making for passive BCI [4]. In particular, the proposed method detected a driver's mental fatigue, *i.e.*, drowsiness, from the EEG signals. Our proposed method predicted 31.10 incorrect trials from a total of 177 samples on average. Furthermore, accurately detecting a drowsy state is one of the most important MSNN capabilities for practical use. Our proposed model only made 5.38 mistakes out of 35 drowsy trials on average, thus exhibiting the highest precision score.

*4) Seizure:* Finally, the MSNN incorrectly identified 5.35 seizures among 178 total test seizure samples. Furthermore, our proposed network was the fastest for detecting seizures, *i.e.*, it exhibited the shortest latency time (approximately 4.98 sec on average) among various methods. In other words, our proposed method demonstrated the best performance even with the shortest latency time. Additionally, the proposed model
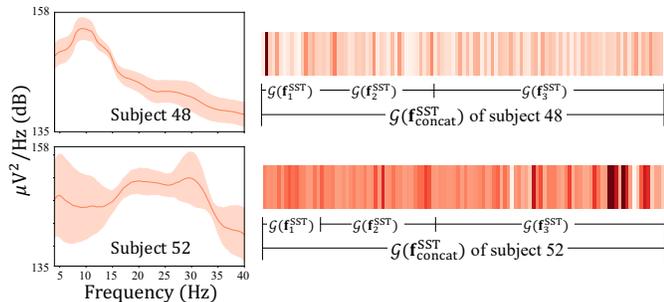
Fig. 3: PSD curves (left) and relevance scores [32] (right) for subject 48 (top) and subject 52 (bottom) from the GIST-MI dataset [8]. For the PSD curves, the solid red line and the shaded region exhibit the mean and standard deviation of PSD values of all trials, respectively. We observed that our proposed MSNN concentrates features from the lower frequency range for subject 48 and a wide range for subject 52.

correctly identified approximately 92% of the seizures within 4.98 sec. We do not present the standard deviation values for this seizure detection experiment because each test trial consisted of different numbers of seizures.

## V. ANALYSES AND DISCUSSIONS

In this section, we analyzed our proposed network. We determined the feature response by estimating PSD values and relevance scores [32] to show the multi-scale learning benefits. We also visualized learned weights and represented features of the proposed method using different methodologies, activation pattern maps [25] and t-SNE plots. Additionally, we observed a practical use for the proposed method, especially for drowsiness and seizure detection experiments.

### A. Multi-Scale EEG Feature Extraction

To demonstrate the multi-scale information capture ability of our proposed method, we estimated and plotted PSD values and relevance scores [32] for MI EEG samples. Specifically, we estimated PSD values for subject 48 and 52 in the GIST-MI dataset [8]'s EEG samples from channels on the motor cortex. Additionally, we calculated relevance scores for those subjects by a *layer-wise relevance propagation* [32]. In our results, all classification methods evenly demonstrated well-generalization (baselines: $\sim$80% and proposed: $\sim$85%) for subject 48, whereas only the proposed method achieved superior performance for subject 52 (baselines: <65% and proposed: $\sim$80%). As Fig. 3 shows, subject 48's EEG samples are highly activated at the $\mu$ range, while subject 52's samples do not show any clear trend at the $\mu$ range, but in a wider range. Our proposed network exhibited a high relevance score at the low-frequency range for subject 48 who exhibited a clear trend at the low-frequency range. Furthermore, the relevance scores for subject 52 were roughly alike for the wider range, where subject 52's PSD demonstrated a less clearly defined trend.

From this phenomenon, we can conclude that our proposed MSNN can capture important features on the multi-scale range, not only in the frequency of interest. In other words, while other existing methods gather spatio-spectral-temporal information at the sequential level, the proposed network exploit multi-scale features, thereby improving learning ability[9].

### B. Activation Patterns

Earlier, Haufe *et al.* [25] proposed an *activation pattern* which is based on a *forward-backward modeling* in signal processing. The activation pattern method [25] provides a way to interpret weight matrices in multivariate neuroimaging, as presented in the signal processing literature.

The proposed method, clearly, decodes the input EEG signal to the corresponding label, *i.e.*, inferring a user's intention or condition from an observed EEG pattern. Therefore, it is a backward process computational model. Hence, for a concrete and meaningful understanding of learned layers, it is essential to reverse this backward process model to a forward process. Finally, in this work, we estimated and visualized the activation patterns of the learned weights shown in Fig. 4a. We extracted the spatial convolutions of Shallow ConvNet [2], Deep ConvNet [2], RSTNN [15], EEGNet [3], and the proposed model. Then, we estimated activation patterns and visualized them in a topological manner. We do not estimate ESTCNN [23] activation patterns because the ESTCNN [23] does not have any spatial feature representation layers and those visualized patterns are estimated by the first subject's first fold data in the GIST-MI dataset [8]. Finally, we normalized the activation patterns in [0, 1] range before visualization.
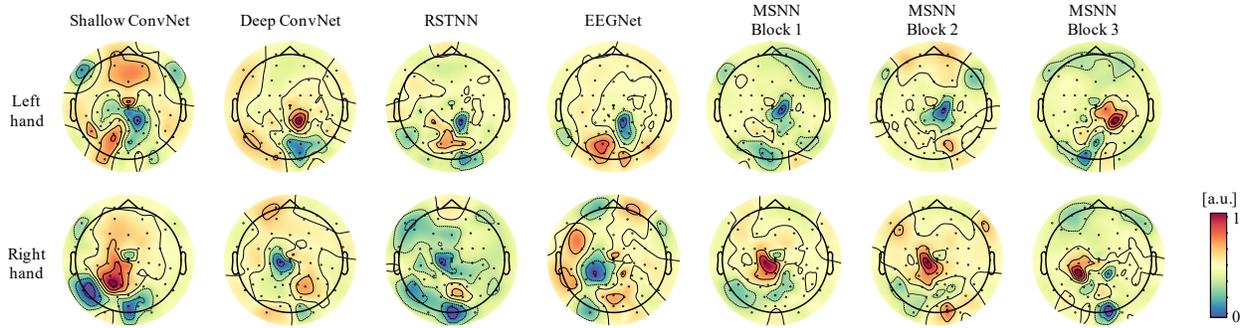
In this investigation, we observed right-lateralized brain activation/deactivation patterns, and the same patterns in the left hemisphere when a user imagined the movement of left-hand and right-hand respectively. Furthermore, the proposed model shows relatively clearer patterns than the other models, thus, we can conclude that our method thoroughly represents input EEG signal spatial features.

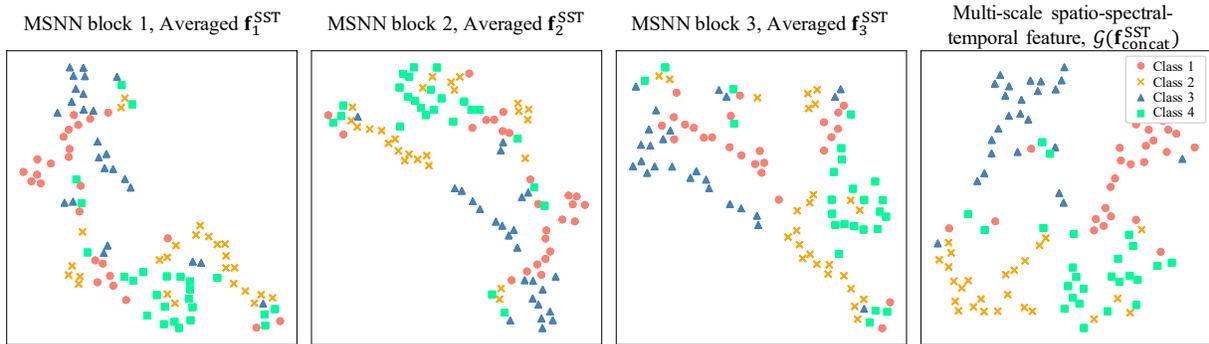### C. Discriminative Power of EEG Representations

To validate the representation ability of the proposed network, we plotted t-SNE transformed learned features shown in Fig. 4b. Specifically, we exhibited extracted features from test SSVEP EEG samples from the first, second, and third spatio-spectral-temporal feature representation layers, *i.e.*, $\mathbf{f}_1^{SST}$, $\mathbf{f}_2^{SST}$, and $\mathbf{f}_3^{SST}$ (first three figures in Fig. 4b). Then, we also depicted the final learned feature, *i.e.*, $\mathcal{G}(\mathbf{f}_{concat}^{SST})$. These intermediate features $\mathbf{f}_1^{SST}$, $\mathbf{f}_2^{SST}$, and $\mathbf{f}_3^{SST}$ are temporally pooled just for visualization like $\mathcal{G}(\mathbf{f}_{concat}^{SST})$. We used the first subject's first session data in the KU-SSVEP dataset [6], and used a learning rate of 200, a perplexity of 10 for the t-SNE calculation, and visualization.

From these visualized represented features, we could observe that $\mathcal{G}(\mathbf{f}_{concat}^{SST})$ is more class-discriminative than the other intermediate features. Additionally, we observed a trend, which demonstrated that a feature learned by a deeper layer is more disentangled than others learned by shallower layers.
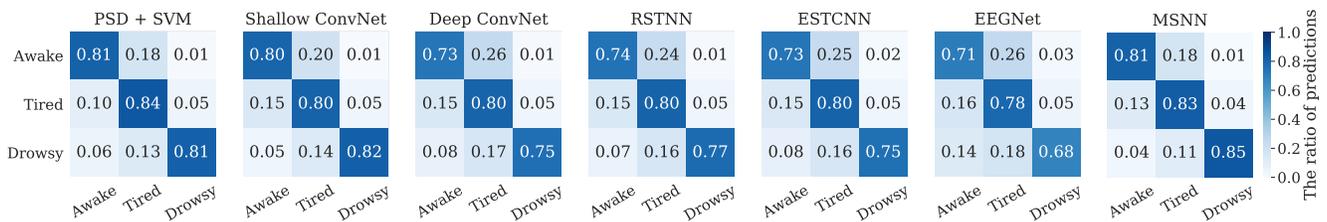
[9]Randomly selected additional results are reported in Supplementary C.
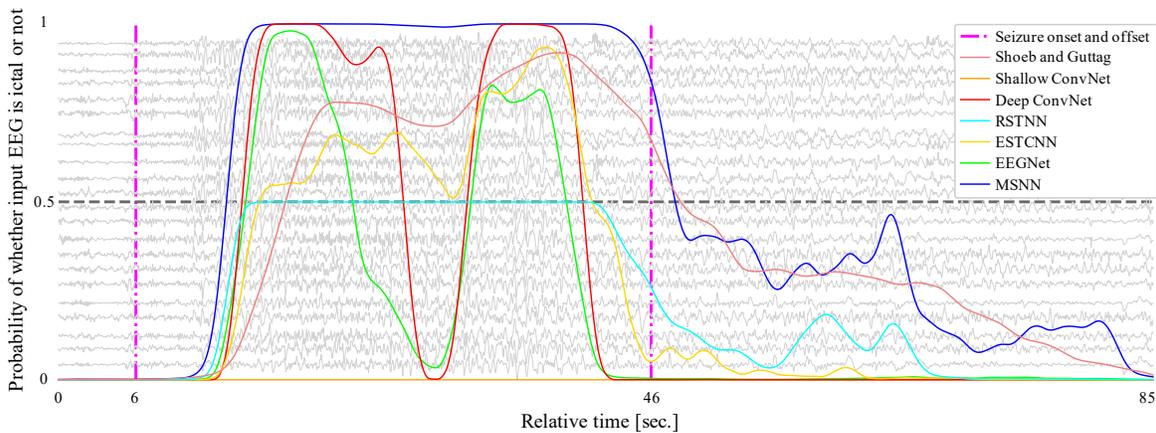
(a) Topologically visualized activation pattern maps [25] of comparable baselines, and three spatial convolutions in the proposed network. All these visualized patterns here are estimated by the first subject's first fold EEG signals in the GIST-MI dataset [8] and normalized in a range between 0 and 1. Finally, [a.u.] denotes an arbitrary unit.



(b) Visualization of t-SNE transformed represented features for test SSVEP EEG samples. The first three figures denote extracted features by the first, second, and final spatial convolutional layers of the proposed method. The final figure exhibits the GAP [29]-ed feature, $\mathcal{G}(\mathbf{f}_{\text{concat}}^{\text{SST}})$ which is used for final decision making.



(c) Normalized averaged confusion matrices estimated by comparable baselines and the proposed method using the SEED-VIG dataset [10].



(d) Changes of probabilities estimated by comparable baselines, and the proposed method. These plots demonstrate the probability of whether input EEG is ictal or not. Two dot-dashed lines (magenta) denote the seizure onset and ending, respectively, labeled by Shoeb [11].

Fig. 4: Investigation of learned weights (Fig. 4a) and represented features (Fig. 4b), and inspection of the practical usage of the proposed network (Fig. 4d and 4c).

## D. Mental Fatigue Classification

For the application analysis of drowsiness detection, we visualized confusion matrices that were estimated by the experimental results of the SEED-VIG dataset [10] in Fig. 4c. Because the labels that identify the mental status were decided using the PERCLOS levels [10], the label at the boundary of the two classes may not be accurate. In this respect, we can conclude that the proposed method is useful for drowsiness state detection because false detections predicted by the proposed method are mostly at the boundaries between classes, *e.g.*, the 'awake' *vs.* 'tired' or 'tired' *vs.* 'drowsy' case. In addition, for practical application, it is essential to detect the drowsy state accurately to avoid unexpected situations, such as a car accident. The proposed method achieved the highest and most promising result for detecting drowsiness among other baselines, *i.e.*, it achieved the highest precision score for identifying the drowsy state. Therefore, we can also expect that our proposed method can be applied in real-world situations.

## E. Early Seizure Detection

Early detection [27] of seizures is one of the most important potential practical applications for this work. Hence, we also validated tthe benefits of the proposed method in early seizure detection. Specifically, in the training phase, the MSNN was trained using normal and ictal EEG samples with binary labels (*e.g.*, 0: normal and 1: seizure) similar to a conventional training framework. In the testing phase, we input the EEG samples using a sliding window with a 1/256 stride. Then, we observed the change in the output probability values to determine the character of the input (normal or ictal).

Additionally, we visualized these changes in Fig. 4d (We used the first subject's third EEG trial in the CHB-MIT dataset [11] for the visualization). In Fig. 4d, magenta-colored dot-dashed lines denote the seizure onset and offset. Colored solid lines denote the probability change of various methods. In this visualization, we observed that the proposed method is more stable for detecting seizures. Specifically, the proposed method detects the seizure EEG signal as a seizure state with a strong probability (almost 1), whereas the other methods have low confidence values (Shoeb and Guttag [14]'s method and ESTCNN [23]) or even make incorrect decisions regarding the seizure state (Shallow ConvNet [2], Deep ConvNet [2], RSTNN [15], and EEGNet [3]).

## VI. Conclusion

In this work, we proposed a novel and compact deep multi-scale neural network which can learn multi-scale EEG signal features. In our experiments, we validated our novel architecture's effectiveness over diverse EEG paradigms, MI, SSVEP, seizure, and drowsy EEG signals. Furthermore, we inspected the relevance scores to demonstrate the benefits of the multi-scale feature extraction ability, investigated activation pattern maps to understand what types of neurophysiological phenomena were learned by our CNN model, and visualized the t-SNE of learned features to examine the ability of our method to differentiate feature classes. Finally, we also

demonstrated that the proposed method can be used for precise drowsiness detection and early seizure detection. In all these respects, we concluded that the proposed deep multi-scale neural network offers significant potential for interpreting EEG signals. Additionally, because the proposed network is clearly *generalizable* to various EEG paradigms, it is expected to have promising benefits that can apply to neural architecture search methods [33], thereby making a deep learning-based BCI adaptable to different paradigms.

From a practical standpoint, many limitations remain with regard to the inter-subject variation [24] in performance. In the present work, we experimented in a subject-dependent manner. In general use, it is important for a BCI system to be useful for any subject operating in a subject-independent way. Thus, in the future, we will focus on developing a subject-neutral multi-paradigm BCI system using adversarial learning [34], [35] or other learning strategies [36].

## References

[1] B. Blankertz, R. Tomioka, S. Lemm, M. Kawanabe, and K.-R. Muller, "Optimizing Spatial Filters for Robust EEG Single-trial Analysis," *IEEE Signal Processing Magazine*, vol. 25, no. 1, pp. 41–56, 2008.

[2] R. T. Schirrmeister, J. T. Springenberg, L. D. J. Fiederer, M. Glasstetter, K. Eggensperger, M. Tangermann, F. Hutter, W. Burgard, and T. Ball, "Deep Learning with Convolutional Neural Networks for EEG Decoding and Visualization," *Human Brain Mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.

[3] V. J. Lawhern, A. J. Solon, N. R. Waytowich, S. M. Gordon, C. P. Hung, and B. J. Lance, "EEGNet: A Compact Convolutional Neural Network for EEG-based Brain–Computer Interfaces," *Journal of Neural Engineering*, vol. 15, no. 5, p. 056013, 2018.

[4] P. Aricò, G. Borghini, G. Di Flumeri, N. Sciaraffa, and F. Babiloni, "Passive BCI Beyond the Lab: Current Trends and Future Directions," *Physiological Measurement*, vol. 39, no. 8, p. 08TR02, 2018.

[5] X. Zhang, L. Yao, X. Wang, J. Monaghan, and D. Mcalpine, "A Survey on Deep Learning based Brain Computer Interface: Recent Advances and New Frontiers," *arXiv preprint arXiv:1905.04149*, 2019.

[6] M.-H. Lee, O.-Y. Kwon, Y.-J. Kim, H.-K. Kim, Y.-E. Lee, J. Williamson, S. Fazli, and S.-W. Lee, "EEG Dataset and OpenBMI Toolbox for Three BCI Paradigms: An Investigation into BCI Illiteracy," *GigaScience*, vol. 8, no. 5, p. giz002, 2019.

[7] M. Nakanishi, Y. Wang, Y.-T. Wang, and T.-P. Jung, "A Comparison Study of Canonical Correlation Analysis based Methods for Detecting Steady-State Visual Evoked Potentials," *PLoS One*, vol. 10, no. 10, p. e0140703, 2015.

[8] H. Cho, M. Ahn, S. Ahn, M. Kwon, and S. C. Jun, "EEG Datasets for Motor Imagery Brain–Computer Interface," *GigaScience*, vol. 6, no. 7, p. gix034, 2017.

[9] C. O'Reilly, N. Gosselin, J. Carrier, and T. Nielsen, "Montreal Archive of Sleep Studies: An Open-access Resource for Instrument Benchmarking and Exploratory Research," *Journal of Sleep Research*, vol. 23, no. 6, pp. 628–635, 2014.

[10] W.-L. Zheng and B.-L. Lu, "A Multimodal Approach to Estimating Vigilance using EEG and Forehead EOG," *Journal of Neural Engineering*, vol. 14, no. 2, p. 026017, 2017.

[11] A. H. Shoeb, "Application of Machine Learning to Epileptic Seizure Onset Detection and Treatment," Ph.D. dissertation, Massachusetts Institute of Technology, 2009.

[12] H.-I. Suk and S.-W. Lee, "A Novel Bayesian Framework for Discriminative Feature Extraction in Brain-Computer Interfaces," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 2, pp. 286–299, 2012.

[13] K. K. Ang, Z. Y. Chin, H. Zhang, and C. Guan, "Filter Bank Common Spatial Pattern (FBCSP) in Brain-Computer Interface," in *IEEE International Joint Conference on Neural Networks*, 2008, pp. 2390–2397.

[14] A. H. Shoeb and J. V. Guttag, "Application of Machine Learning to Epileptic Seizure Detection," in *Proceedings of the 27th International Conference on Machine Learning*, 2010, pp. 975–982.

[15] W. Ko, J. Yoon, E. Kang, E. Jun, J.-S. Choi, and H.-I. Suk, "Deep Recurrent Spatio-Temporal Neural Network for Motor Imagery based BCI," in *2018 6th International Conference on Brain-Computer Interface (BCI)*, 2018, pp. 1–3.

[16] F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1251–1258.

[17] A. Supratak, H. Dong, C. Wu, and Y. Guo, "DeepSleepNet: A Model for Automatic Sleep Stage Scoring based on Raw Single-channel EEG," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 11, pp. 1998–2008, 2017.

[18] S. Sakhavi, C. Guan, and S. Yan, "Learning Temporal Information for Brain-Computer Interface Using Convolutional Neural Networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2018.

[19] N. Waytowich, V. J. Lawhern, J. O. Garcia, J. Cummings, J. Faller, P. Sajda, and J. M. Vettel, "Compact Convolutional Neural Networks for Classification of Asynchronous Steady-State Visual Evoked Potentials," *Journal of Neural Engineering*, vol. 15, no. 6, p. 066031, 2018.

[20] N.-S. Kwak, K.-R. Müller, and S.-W. Lee, "A Convolutional Neural Network for Steady State Visual Evoked Potential Classification Under Ambulatory Environment," *PLoS one*, vol. 12, no. 2, p. e0172578, 2017.

[21] U. Asif, S. Roy, J. Tang, and S. Harrer, "SeizureNet: A Deep Convolutional Neural Network for Accurate Seizure Type Classification and Seizure Detection," *arXiv preprint arXiv:1903.03232*, 2019.

[22] A. Emami, N. Kunii, T. Matsuo, T. Shinozaki, K. Kawai, and H. Takahashi, "Seizure Detection by Convolutional Neural Network-based Analysis of Scalp Electroencephalography Plot Images," *NeuroImage: Clinical*, vol. 22, p. 101684, 2019.

[23] Z. Gao, X. Wang, Y. Yang, C. Mu, Q. Cai, W. Dang, and S. Zuo, "EEG-Based Spatio-Temporal Convolutional Neural Network for Driver Fatigue Evaluation," *IEEE Transactions on Neural Networks and Learning Systems*, 2019.

[24] V. Jayaram, M. Alamgir, Y. Altun, B. Scholkopf, and M. Grosse-Wentrup, "Transfer Learning in Brain-Computer Interfaces," *IEEE Computational Intelligence Magazine*, vol. 11, no. 1, pp. 20–31, 2016.

[25] S. Haufe, F. Meinecke, K. Görgen, S. Dähne, J.-D. Haynes, B. Blankertz, and F. Bießmann, "On the Interpretation of Weight Vectors of Linear Models in Multivariate Neuroimaging," *NeuroImage*, vol. 87, pp. 96–110, 2014.

[26] T. H. Sanders, M. McCurry, and M. A. Clements, "Sleep Stage Classification with Cross Frequency Coupling," in *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, 2014, pp. 4579–4582.

[27] J. Lee, J. Park, S. Yang, H. Kim, Y. S. Choi, H. J. Kim, H. W. Lee, and B.-U. Lee, "Early Seizure Detection by Applying Frequency-based Algorithm Derived from the Principal Component Analysis," *Frontiers in Neuroinformatics*, vol. 11, p. 52, 2017.

[28] M. Liang and X. Hu, "Recurrent Convolutional Neural Network for Object Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3367–3375.

[29] M. Lin, Q. Chen, and S. Yan, "Network in Network," *arXiv preprint arXiv:1312.4400*, 2013.

[30] X. Glorot and Y. Bengio, "Understanding the Difficulty of Training Deep Feedforward Neural Networks," in *Proceedings of the thirteenth International Conference on aAtificial Intelligence and Statistics*, 2010, pp. 249–256.

[31] A. Binder, S. Bach, G. Montavon, K.-R. Müller, and W. Samek, "Layer-Wise Relevance Propagation for Deep Neural Network Architectures," in *Information Science and Applications*. Springer, 2016, pp. 913–922.

[32] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller, "Explaining Nonlinear Classification Decisions with Deep Taylor Decomposition," *Pattern Recognition*, vol. 65, pp. 211–222, 2017.

[33] E. Rapaport, O. Shriki, and R. Puzis, "EEGNAS: Neural Architecture Search for Electroencephalography Data Analysis and Decoding," in *International Workshop on Human Brain and Artificial Intelligence*. Springer, 2019, pp. 3–20.

[34] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-Adversarial Training of Neural Networks," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.

[35] E. Jeon, W. Ko, and H.-I. Suk, "Domain Adaptation with Source Selection for Motor-Imagery based BCI," in *2019 7th International Winter Conference on Brain-Computer Interface (BCI)*, 2019, pp. 1–4.

[36] Z. Li and D. Hoiem, "Learning without Forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.