

# Functional Neural Networks: Shift invariant models for functional data with applications to EEG classification

Florian Heinrichs<sup>1</sup> Mavin Heim<sup>1</sup> Corinna Weber<sup>1</sup>

## Abstract

It is desirable for statistical models to detect signals of interest independently of their position. If the data is generated by some smooth process, this additional structure should be taken into account. We introduce a new class of neural networks that are shift invariant and preserve smoothness of the data: functional neural networks (FNNs). For this, we use methods from functional data analysis (FDA) to extend multi-layer perceptrons and convolutional neural networks to functional data. We propose different model architectures, show that the models outperform a benchmark model from FDA in terms of accuracy and successfully use FNNs to classify electroencephalography (EEG) data.

## 1. Introduction

When autonomous machines act in their environment or humans interact with computers, the necessary input data is often streamed continuously. In some settings, the input streams can be easily transformed into control signals based on simple physical models. However, in more advanced scenarios, it is necessary to develop a more complex, data-driven model and use its predictions to control the machine. An important example of the latter scenario are brain-computer interfaces (BCIs), which record the user's brain activity and decode control signals based on the measured data.

Most statistical models require an input of fixed dimension and a common approach is to extract windows of a fixed size with a fixed step size from the continuous data stream. These *sliding windows* are then used to predict the desired control signal with one classification per window. The advantage of this approach is that only the most recent data is used for predictions, but it comes at cost: the signal of interest might

occur at any time point in the extracted window. Thus, any classifier of the sliding windows needs to be “shift invariant” in the sense that it detects the desired signal independently of its position in the window.

In the context of BCIs, more specifically for the analysis of data measured via electroencephalography (EEG), traditional methods are based on carefully selected features that are calculated from the data. Commonly applied techniques include the principal component analysis of EEG data in the time domain, or features based on the power spectrum in the frequency domain (Azlan & Low, 2014; Boubchir et al., 2017; Boonyakitanont et al., 2020). Due to recent advances in the field of deep learning, different architectures of neural networks were proposed that avoid a manual feature extraction and seem to outperform more traditional methods. For example the neural network *EEGNet* was proposed to support multiple BCI paradigms and is often referred to as benchmark model in the field (Lawhern et al., 2018). In a clinical setting, some variant of the VGG16 neural network was used to detect signals associated with epilepsy (da Silva Lourenço et al., 2021). In general, deep learning has been applied successfully to a variety of tasks related to EEG data (Craik et al., 2019; Roy et al., 2019).

Inspired by their successes in computer vision and natural language processing, common neural networks used for the classification of EEG data are based on convolutions. Convolutional neural networks are often considered to be shift invariant and are therefore a good choice in the given context. However, they do not take the specific structure of EEG data into account. Similar to most physical processes, the electrical activity, that is recorded on the user's scalp by the EEG, can be considered as smooth (Ramsay & Silverman, 2005). To reflect this additional structure, it is more appropriate to model the data as a (discretized) sample of an underlying smooth function.

The latter paradigm is the basis of *functional data analysis* (FDA), a branch of statistics that received more and more attention throughout the last decades and remains an active area of research. Most concepts from multivariate statistics have been extended to functional data (Ramsay & Silverman, 2005; Kokoszka & Reimherr, 2017). For example, many functional versions of the principal component analy-

<sup>1</sup>SNAP GmbH, Gesundheitscampus-Süd 17, 44801 Bochum, Germany. Correspondence to: Florian Heinrichs <mail@florian-heinrichs.de>, f.heinrichs@snap-gmbh.com>, Mavin Heim <m.heim@snap-gmbh.com>, Corinna Weber <weber@snap-gmbh.com>.

sis have been proposed in literature (Shang, 2014). Different generalizations of the linear model to functional covariates and/or functional responses have been introduced (Cardot et al., 1999; Cuevas et al., 2002). Finally, Portmanteau-type tests for detecting serial correlation have been proposed for functional time series (Gabrys & Kokoszka, 2007; Bücher et al., 2020). This functional approach allows it to extract previously unavailable information from the data in form of derivatives of the continuous signal.

As methods from FDA take the functional structure of physical processes into account, they would be suitable classifiers. However, classic methods from FDA are in general not shift invariant and require the signal of interest to be at a fixed point in time. In some applications it is possible to “register” the functions through a suitable transformation of time (Sakoe & Chiba, 1978; Kneip & Gasser, 1992; Gasser & Kneip, 1995; Ramsay & Li, 1998). However, in the context of sliding windows curve registration is often not feasible and methodology that requires previous registration cannot be applied reliably.

In the present work, we propose a framework that combines the advantages of neural networks (particularly CNNs) and FDA: functional neural networks (FNNs). On one side, these networks are shift invariant, and on the other side, they are able to model the functional structure of their input. FNNs have several advantages over scalar-valued neural networks. They are independent of the sample frequency of the input data, as long as the input can be rescaled to a certain interval. Further, they allow to predict smooth outputs. And finally, they are more transparent to some extent due to smoothness constraints.

We summarize our contribution as follows:

- We propose extensions of fully-connected and convolutional layers to functional data.
- We present architectures of functional neural networks based on these extensions.
- We show that the proposed methodology works through a simulation study and real data experiments.

Whereas multi-layer perceptrons (MLPs) are not shift invariant, the introduced functional convolutional layers allow the construction of shift invariant functional convolutional neural networks. This makes FNNs helpful in any scenario where sliding windows based on a (possibly multivariate) continuous data stream are classified, and they can be employed in a variety of applications.

## 2. Related Work

To the best of our knowledge, the combination of functional data and (convolutional) neural networks is only discussed

in a handful of papers and the proposed methodology extends previous results. In early works MLPs with functional inputs and neurons, that transform the functional data to scalar values in the first layer, were introduced (Rossi et al., 2002; Rossi & Conan-Guez, 2005; Rossi et al., 2005). (Zhao, 2012) proposed an algorithm to train similar MLPs with inputs from a real Hilbert space. Subsequently, (Wang et al., 2019) proposed to use functional principal components for the transformation of the functional inputs to scalar values in the first layer. (Wang et al., 2020) added another layer based on functional principal components to transform the scalar-valued output of the MLP back to functional data in the last layer. More recently, fully functional neurons were proposed (Rao & Reimherr, 2021b;a).

## 3. Mathematical Preliminaries

Let us assume, we observe  $d$  quantities at  $T$  time instants for  $N \in \mathbb{N}$  individuals, providing us with matrices of observations

$$\mathbf{X}^{(n)} = \begin{pmatrix} X_{1,1}^{(n)} & \cdots & X_{1,T}^{(n)} \\ \vdots & \ddots & \vdots \\ X_{d,1}^{(n)} & \cdots & X_{d,T}^{(n)} \end{pmatrix},$$

for  $n = 1, \dots, N$ , and jointly with  $\mathbf{X}^{(n)}$  their corresponding “labels”  $\mathbf{Y}^{(n)}$  which might be vectors in  $\mathbb{R}^c$  or matrices

$$\mathbf{Y}^{(n)} = \begin{pmatrix} Y_{1,1}^{(n)} & \cdots & Y_{1,T}^{(n)} \\ \vdots & \ddots & \vdots \\ Y_{c,1}^{(n)} & \cdots & Y_{c,T}^{(n)} \end{pmatrix}$$

where  $c$  denotes the number of quantities that we observe for  $\mathbf{Y}^{(n)}$ . Further assume the observed quantities to be noisy versions of an underlying smooth signal, i.e.,

$$X_{i,t}^{(n)} = f_i^{(n)}\left(\frac{t}{T}\right) + \varepsilon_{i,t}^{(n)}, \quad (1)$$

for smooth functions  $f_i^{(n)}$  and centered errors  $\varepsilon_{i,t}^{(n)}$ , for  $i = 1, \dots, d, t = 1, \dots, T$  and  $n = 1, \dots, N$ . Note that the degree of smoothness might vary for different applications, which leads to slight modifications in the model. This representation suggests the use of methods from functional data analysis, which consider the intrinsic structure of the data.

Throughout this work, we only require the functions  $f_i^{(n)}$  to be square-integrable, i.e.,  $f_i^{(n)} \in L^2([0, 1]) = \{f : [0, 1] \rightarrow \mathbb{R} : \int_0^1 f(x)^2 dx < \infty\}$ . Similarly, in case of matrix-valued labels  $\mathbf{Y}^{(n)}$ , we assume their entries to be discretized versions of some underlying functions  $g_i^{(n)} \in L^2([0, 1])$ , more specifically  $Y_{i,t}^{(n)} = g_i^{(n)}(t/T)$ .

Our aim is to approximate the functional  $F : (L^2([0, 1]))^d \rightarrow \mathcal{Y}$ , which maps an observation  $\mathbf{X}$  to its

corresponding label  $\mathbf{Y}$ , where  $\mathcal{Y} = \mathbb{R}^c$  for vector-valued labels  $\mathbf{Y}^{(n)}$  and  $\mathcal{Y} = (L^2([0, 1]))^c$  for matrix-valued labels. Formally, the functional  $F$  corresponds to the conditional expectation  $\mathbb{E}[\mathbf{Y}|\mathbf{X}]$ .

In case of classification problems, each coordinate  $F_i(\mathbf{X})$  of the functional  $F$  can be interpreted as the probability of  $\mathbf{X}$  belonging to class  $i \in \{1, \dots, c\}$ .

### 3.1. Preprocessing

Before putting observations into a neural network, it is often helpful to preprocess them by applying certain filters or normalization. In our case, we work with noisy functional data, observed at discrete time points. In functional data analysis, a common first step for this kind of data is *smoothing*, which helps to reduce the errors and extends the observations from discrete time points to a continuous interval. Another preprocessing step frequently used for neural networks, is some form of normalization to ensure that the data is of a similar magnitude. We employ *local linear estimation* for smoothing the data and *standardization* for its normalization as described below.

#### 3.1.1. SMOOTHING

In the literature, there exists a variety of smoothing procedures from Fourier series to expansions based on B-splines or wavelets (Ramsay & Silverman, 2005). We use local polynomial regression to estimate the functions  $f_i^{(n)}$  and their first derivative(s) (Fan & Gijbels, 1996).

For the sake of clarity, we omit some indices and rewrite (1) as  $X_t = f(\frac{t}{T}) + \varepsilon_t$  for a moment. Then, if  $f$  is  $p+1$  times differentiable with bounded derivatives, we can define the local polynomial estimator as

$$\begin{aligned} & (\hat{f}(x), \hat{f}'(x), \dots, \widehat{f^{(p)}}(x)) \\ &= \arg \min_{\beta_0, \dots, \beta_p} \sum_{t=1}^T \left( X_t - \sum_{j=0}^p \beta_j \left( \frac{t}{T} - x \right)^j \right)^2 K_h \left( \frac{t}{T} - x \right) \end{aligned}$$

to estimate  $f$  and its first  $p$  derivatives. Here  $K$  denotes a kernel function,  $h$  the bandwidth of the estimator and  $K_h(\cdot) = K(\frac{\cdot}{h})$ . In the following, we assume  $K : \mathbb{R} \rightarrow \mathbb{R}$  to be a symmetric, twice differentiable function, supported on the interval  $[-1, 1]$  and satisfying  $\int_{[-1, 1]} K(x) dx = 1$ .

From the above definition, explicit formulas can be derived for the estimators by calculating the derivatives of the right-hand side with respect to  $\beta_j$  ( $j = 1, \dots, p$ ) and a Taylor expansion. The result can be rewritten as a convolution of the signal and a filter depending on the kernel  $K$  and the bandwidth  $h$ .

To simplify the notation, we will refer to the estimators of the functions  $f_i^{(n)}$  and their derivatives as  $h_{i,1}^{(n)}$ , thus,

we obtain estimators  $(h_{i,1}^{(n)}, (h_{i,1}^{(n)})', \dots, (h_{i,1}^{(n)})^{(p)})$  for each  $f_i^{(n)}$ ,  $i = 1, \dots, d$ ,  $n = 1, \dots, N$ .

The choice of the bandwidth is crucial in order to obtain a good estimate of the underlying functions. If the bandwidth is chosen too small, the estimator will overfit the data, whereas a large bandwidth leads to over-smoothing (Silverman, 2018). Oftentimes it is a good idea to use cross validation to select a bandwidth that minimizes a certain error measure, such as the mean squared error. Generally, the estimation of higher derivatives requires larger bandwidths than the estimation of the function itself.

#### 3.1.2. NORMALIZATION

When neural networks are trained via some form of gradient descent, it is crucial to ensure that the input data is of a similar size, which is done through prior normalization. There are many different normalization methods and the most useful choice depends on the specific application. In the following, we will standardize the data by subtracting the mean and dividing by the standard deviation across a suitable range of the data. As we did not make any assumptions about the relation between the signals  $f_i^{(n)}$  and  $f_j^{(n)}$ , we standardize each smoothed signal  $h_{i,1}^{(n)}$  (and its derivatives) separately, i. e., we calculate

$$h_{i,2}^{(n)} = \frac{h_{i,1}^{(n)} - \int_0^1 h_{i,1}^{(n)}(x) dx}{\left( \int_0^1 \left( h_{i,1}^{(n)}(x) - \int_0^1 h_{i,1}^{(n)}(y) dy \right)^2 dx \right)^{1/2}}.$$

After this transformation, the signals are of a similar magnitude, for each observation  $\mathbf{X}^{(n)}$ .

## 4. Functional Layers

### 4.1. Functional Multilayer Perceptrons

Once the data is smoothed and prepared to be analyzed as functional data, it is not clear how to design neural networks that take this additional structure into account. The simplest form of an artificial neural network with scalar input  $(h_1, \dots, h_d)$  is the multilayer perceptron, that consists of  $L$  layers with  $J_1, J_2, \dots, J_L$  neurons each. The value at neuron  $k$  in the  $\ell$ -th layer is then calculated as

$$H_{(k)}^{(\ell)} = \sigma \left( b_{(k)}^{(\ell)} + \sum_{j=1}^{J_{\ell-1}} w_{(j,k)}^{(\ell)} H_{(j)}^{(\ell-1)} \right),$$

where  $H_{(k)}^{(0)} = h_k$  denotes the network's input,  $H_{(k)}^{(L)}$  its output,  $b_{(k)}^{(\ell)}$  the  $k$ th neuron's bias and  $w_{(j,k)}^{(\ell)}$  the weight between the  $k$ th neuron in layer  $\ell$  and the  $j$ th neuron in layer  $\ell - 1$ . The function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is referred to as

activation function and enables the network to reflect non-linear dependencies.

When the input data is not scalar, but functional, (Rao & Reimherr, 2021a) propose to replace the scalar biases by functional biases and the weights between neurons by integral kernels, finally defining the neurons' values as

$$H_{(k)}^{(\ell)}(s) = \sigma \left( b_{(k)}^{(\ell)}(s) + \sum_{j=1}^{J_{\ell-1}} \int w_{(j,k)}^{(\ell)}(s,t) H_{(j)}^{(\ell-1)}(t) dt \right). \quad (2)$$

While this extension allows to model rather general relations between the model's input and its desired output, this flexibility makes the network's training difficult because we need to find optimal weight functions for any connection between two neurons.

Starting from the fully-connected multilayer perceptron, many advances in deep learning are due to more specific architectures, which reduce the number of model parameters to mitigate the curse of dimensionality. For instance, convolutional neural networks can be interpreted as multi-layer perceptrons, where most weights vanish and the remaining connections between neurons share a smaller set of weights.

In a similar fashion, we propose to simplify the neuron model in (2) by using simple weight functions  $w_{(j,k)}^{(\ell)} : [0, 1] \rightarrow \mathbb{R}$  rather than integral kernels in  $L^2([0, 1]^2)$ . This adaptation leads to neurons defined via

$$H_{(k)}^{(\ell)}(t) = \sigma \left( b_{(k)}^{(\ell)}(t) + \sum_{j=1}^{J_{\ell-1}} w_{(j,k)}^{(\ell)}(t) H_{(j)}^{(\ell-1)}(t) \right). \quad (3)$$

The above defined neurons are fully functional in the sense that both their input and output are functions. If we try to predict scalar-valued labels in  $\mathbb{R}^c$ , we need to summarize the information contained in the functions. We propose to calculate the scalar product of the weights and their corresponding inputs, leading to

$$H_{(k)}^{(\ell)} = \sigma \left( b_{(k)}^{(\ell)} + \sum_{j=1}^{J_{\ell-1}} \int w_{(j,k)}^{(\ell)}(t) H_{(j)}^{(\ell-1)}(t) dt \right). \quad (4)$$

With this definition of a functional multilayer perceptron (F-MLP), we simplified the training and need to optimize functional weights of one variable. The theoretical framework to train the model through backpropagation based on Fréchet derivatives is provided by (Rossi et al., 2002; Olver, 2016; Rao & Reimherr, 2021a).

The computation of Fréchet derivatives becomes tedious and computationally expensive. An efficient approach to simplify computations and simultaneously reduce the dimension of the weights' space, is to replace the weights

$w_{(j,k)}^{(\ell)}(t)$  by linear combinations of a finite set of base functions. Therefore, let  $\{\varphi_i\}_{i=0}^q$  be a set of suitable functions, such as Legendre polynomials, wavelets or the first  $q/2$  sine-cosine pairs of the Fourier basis, and consider the linear combination

$$w_{(j,k)}^{(\ell)}(t) = \sum_{i=0}^q w_{(j,k)}^{(\ell,i)} \varphi_i(t), \quad (5)$$

for some scalar weights  $w_{(j,k)}^{(\ell,i)}$ . With this representation, the fully functional neural network can be described through scalar weights and we are able to use the standard scalar backpropagation.

## 4.2. Functional Convolutional Neural Networks

The functional MLP is particularly useful if the input functions are aligned (or can be aligned via a suitable transformation of time) and the signals of interest happen at the same time instants for all measurements. However, under the sliding window paradigm, for high-noise data such as speech or EEG signals, it is not possible (or at least not useful) to previously register the curves, as the signal of interest may occur at any arbitrary time instant. In this case, MLPs are impractical as they would require many parameters to model complex patterns

For scalar input, alternative network architectures have been developed that are shift invariant and therefore capable to detect certain signals independently of their position. One type of neural network that is considered as "translation invariant" are CNNs, which we can extend to functional data as well.

Similarly to (2), we can define a functional convolutional layer by setting  $w_{(j,k)}^{(\ell)}(s,t) = u_{(j,k)}^{(\ell)}(s-t)$  for some filter (or kernel) function  $u_{(j,k)}^{(\ell)} : \mathbb{R} \rightarrow \mathbb{R}$  with support on  $[-b, b]$  and bandwidth  $b \in (0, 1)$ , ultimately leading to

$$H_{(k)}^{(\ell)}(s) = \sigma \left( b_{(k)}^{(\ell)}(s) + \sum_{j=1}^{J_{\ell-1}} \int u_{(j,k)}^{(\ell)}(s-t) H_{(j)}^{(\ell-1)}(t) dt \right),$$

where the functions  $H_{(k)}^{(\ell)}$  are extended to the interval  $[-b, 1+b]$  by defining them as zero outside of the interval  $[0, 1]$ .

These functional convolutional layers are shift invariant in the sense that a filter, which is capable of detecting a certain signal, would detect it independently of its position in the interval  $[0, 1]$ . Once again, we reduce the dimension of the optimization problem by representing the filters as linear combinations of a set of base functions as in (5).



### 4.3. Architecture

With functional versions of fully connected and convolutional layers at hand, we can define arbitrary architectures of functional neural networks (FNNs). Figure 1 displays FNNs with scalar and functional outputs, respectively. In both architectures, the first layer uses a local linear estimator to smooth the input and estimate derivatives of the smoothed signals, while the second layer standardizes the input across each signal. Following are two functional convolutional layers. For the FNN with scalar output, the last layer is a functional fully connected layer, while for the FNN with functional output, the last layer is a third functional convolutional layer.

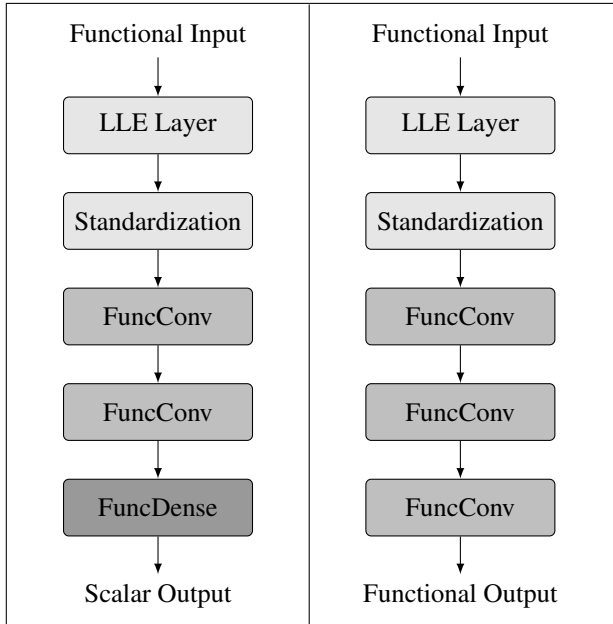


Figure 1. Left: Neural network architecture with functional input and scalar output. Right: Neural network architecture with functional input and output.

## 5. Empirical Results

We now show that the proposed methodology works with simulated and real data and compare it to benchmark models.

### 5.1. Simulation Study

#### SETUP

Inspired by brain activity measured through electroencephalography (EEG), we generate two data sets for the simulation study. In both cases, we simulate two-dimensional samples that belong to one of three classes.

For data set (I), we draw independently frequencies  $\alpha_n \sim$

$\mathcal{U}_{[8,12]}$ ,  $\beta_n \sim \mathcal{U}_{[13,30]}$ , time shifts  $t_1^{(n)}, t_2^{(n)} \sim \mathcal{U}_{[0,1]}$ , class labels  $c_n \sim \mathcal{U}_{\{1,2,3\}}$  and standard normally distributed errors  $\varepsilon_{i,t}^{(n)} \sim \mathcal{N}(0, 1)$ . Based on these random quantities, we construct the continuous signals

$$f_i^{(n)}(x) = (1 - \gamma_i(c_n)) \cdot \sin(2\pi\alpha_n(x + t_i^{(n)})) + \gamma_i(c_n) \cdot \sin(2\pi\beta_n(x + t_i^{(n)}))$$

with class dependent coefficients  $\gamma(1) = (0, 0)$ ,  $\gamma(2) = (0.8, 0.4)$  and  $\gamma(3) = (0.4, 0.8)$  and finally define the discretized, noisy samples  $X_{i,t}^{(n)} = f_i^{(n)}(t/T) + \varepsilon_{i,t}^{(n)}$ . Examples of each class are displayed in Figure 6 of Appendix A.

For data set (II), we draw independently scaling factors  $w_n \sim \mathcal{U}_{[0.05,0.1]}$ , time points  $t_n \sim \mathcal{U}_{[0,1]}$ , class labels  $c_n \sim \mathcal{U}_{\{1,2,3\}}$  and standard normally distributed errors  $\varepsilon_{i,t}^{(n)} \sim \mathcal{N}(0, 1)$ . Based on the scaling factors  $w_n$  and time points  $t_n$ , we construct continuous signals

$$f^{(n)}(x) = \max \left\{ -\frac{4}{w_n^2}(x - t_n)^2 + 3, 0 \right\},$$

which resemble spikes, as displayed in Figure 7 of Appendix A. Again, we define discretized, noisy samples  $X_{i,t}^{(n)} = \gamma_i(c) \cdot f^{(n)}(t/T) + \varepsilon_{i,t}^{(n)}$ , where the class dependent coefficients are  $\gamma(1) = (0, 0)$ ,  $\gamma(2) = (1, 0)$  and  $\gamma(3) = (0, 1)$ .

For both data sets, we vary the sample size  $N \in \{1000, 2000, 3000, 4000, 5000\}$ , while keeping  $T = 250$  fixed.

As baseline models, we use  $k$ -nearest neighbors (KNN) for a varying number of neighbors  $k \in \{1, 2, \dots, 19\}$ . Before feeding the data into the model, we smooth it by applying a local linear estimator as described in Section 3.1.1.

To show that functional neural networks work, and indeed surpass the performance of the baseline models, we use an FNN as described in Section 4.3 with two functional convolutional and one functional dense layer. For the local linear estimation, we use the quartic kernel  $K(x) = \frac{15}{16}(1 - x^2)^2$  with support  $[-1, 1]$  and the bandwidth  $h = 5$  for the estimation of the smooth function and  $h = 10$  for the estimation of its derivative. For each functional layer, we used the first 5 Legendre polynomials as base functions, i.e.  $\varphi_0(x) = 1$ ,  $\varphi_1(x) = x$ ,  $\varphi_2(x) = \frac{1}{2}(3x^2 - 1)$ ,  $\varphi_3(x) = \frac{1}{2}(5x^3 - 3x)$  and  $\varphi_4(x) = \frac{1}{8}(35x^4 - 30x^2 + 3)$ . Further, we used 20 and 10 filters of size 25 for the two convolutional layers. As activation function we chose the *exponential linear unit* (ELU), which is defined as  $\sigma(x) = x \cdot \mathbb{1}(x \geq 0) + (\exp(x) - 1) \cdot \mathbb{1}(x < 0)$ . As loss function, we used the categorical crossentropy.

We trained each model 100 times, while generating a new data set for each trial. The FNNs were trained with 5 epochs.

## RESULTS

The results of the benchmark model for both data sets with a varying number of samples  $N$  and neighbors are displayed in Figures 2 and 3. In Table 1, the results of the benchmark model with the best choice of neighbors  $k$  are compared with the results of the FNN. In all cases, the classifications of the functional neural network are more reliable than those of the  $k$ -nearest neighbors classifier. The FNN achieved an accuracy above 99.6% in all cases, whereas the KNN classifier achieved between 93.0% and 99.3% for data set (I) and between 76.9% and 86.4% for data set (II). As expected, the shift invariance of the FNN makes it particularly helpful for data set (II), where the signal of interest may occur at any point in the observed interval. More specifically, for data set (II) the FNN’s accuracy is at least 13.5% higher than the corresponding accuracy of the KNN classifier.

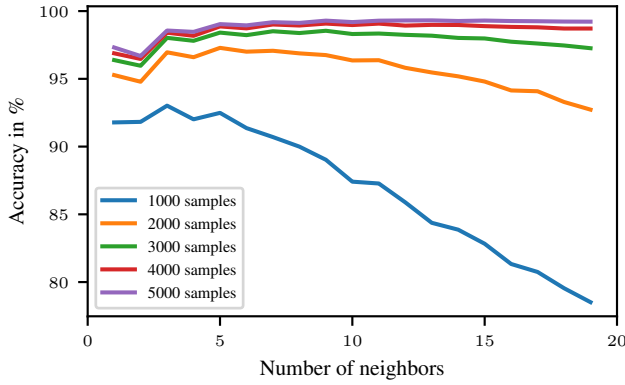


Figure 2. Mean accuracy in percent (y-axis) of the  $k$ -nearest neighbors classifiers for data set (I) and a varying number of neighbors (x-axis).

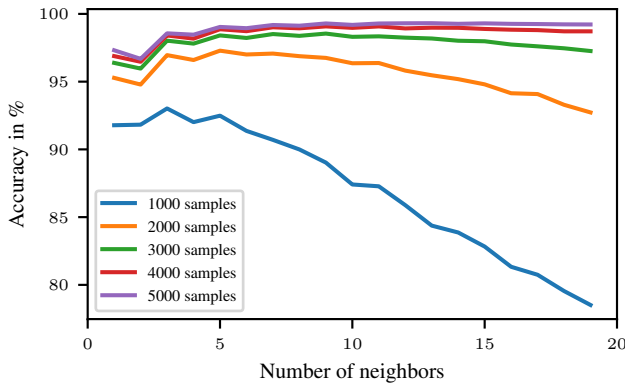


Figure 3. Mean accuracy in percent (y-axis) of the  $k$ -nearest neighbors classifiers for data set (II) and a varying number of neighbors (x-axis).

Table 1. Mean accuracy of the classifiers for the respective data sets.

$N$	DATA SET (I)		DATA SET (II)	
	KNN	FNN	KNN	FNN
1000	93.0%	<b>99.6%</b>	76.9%	<b>99.6%</b>
2000	97.3%	<b>99.8%</b>	81.4%	<b>99.8%</b>
3000	98.5%	<b>99.8%</b>	83.6%	<b>99.8%</b>
4000	99.1%	<b>99.7%</b>	85.5%	<b>99.9%</b>
5000	99.3%	<b>99.8%</b>	86.4%	<b>99.9%</b>

## 5.2. Real Data Experiments

## SETUP

The BCI Competition IV Dataset 2A (Tangemann et al., 2012) is a common benchmark for evaluating the performance of a new method for the analysis of EEG data. To test our method, we used the 9 openly available, labeled recordings of approximately 45 minutes each.

According to the documentation, participants were asked to imagine movements of their left hand (class 1), right hand (class 2), both feet (class 3) and tongue (class 4). During each session, every imaginary movement was repeated 72 times, yielding a total of 288 trials. Each trial took approximately 8 seconds. At the beginning of each trial ( $t = 0$  s), a short acoustic signal and a fixation cross on a black screen appeared. Two seconds later ( $t = 2$  s), a visual cue appeared to indicate the movement, which should be imagined. The imaginary movement can be assumed to start approximately half a second after the cue ( $t = 2.5$  s) and end when the fixation cross disappeared ( $t = 6$  s). Each trial was followed by a short break to separate it from subsequent trials. The participants’ brain activity was measured through a 22-channel EEG with 3 additional EOG (*Electrooculography*) channels at a sampling rate of 250 Hz.

For this data set, the *classic approach* to benchmark a new method is to cut windows from each trial, e. g., between 2.5 s and 4.5 s after trial onset, which is feasible since the trial and cue onsets are known. However, if we move beyond externally triggered actions, we need another approach. This is particularly important in the case of brain-computer interfaces where devices should be controlled continuously. In this case, a common approach is to use *sliding windows*, i. e., to use overlapping windows of a fixed length with a fixed step size.

We tested the proposed functional neural network, as described in Section 4.3 with the same specifications as in the simulation study, and compared it to the EEGNet with its default choice of hyperparameters as suggested by (Lawhern et al., 2018). However, to account for the different degrees of complexity of the EEG data, we chose to use different numbers of filters in the convolutional layers. For the classic

approach with 4 classes (corresponding to the 4 different imaginary movements), we tested two models with 5 + 10 and 3 + 12 filters, respectively, and denote these models by FNN(5, 10) and FNN(3, 12). These models have 2,344 and 1,564 trainable weights, which is slightly less than the 2,526 trainable parameters of the EEGNet. For the sliding windows approach with 7 classes, we increased the number of filters to 40 + 20 and 5 + 10. In this setting, we get models with 19,767 and 2,497 trainable weights compared to 2,783 for the EEGNet (with 7 classes).

For the classic approach, we used windows between the cue onset ( $t = 2$  s) and the disappearance of the fixation cross ( $t = 6$  s). We split each recording into 80% train and 20% test data, which corresponds to 230 and 58 windows each. We trained the proposed FNN and the EEGNet with 250 batches of size 32 to distinguish between the four classes (left hand, right hand, feet, tongue). In total 8000 samples were used, which means that each of the 230 training windows was used approximately 35 times.

For the sliding windows approach, we split each recording into 80% train and 20% test data, which corresponds to approximately 36 and 9 minutes respectively. We used windows of 1 s and a step size of 0.016 s which led to more than 125,000 sliding windows. These windows might coincide with a break between trials (class 1), the time between trial and cue onset (class 3), the time between cue onset and imagined movement (class 2) or one of the four imagined movements (classes 4 - 7). The windows at the transition between two classes were labeled with the most frequent class. This problem is substantially more complex than the classic approach and we have more varied data. Thus, we trained both models with 4000 batches of size 32, which is close to the number of sliding windows.

We trained each model 10 times for each of the 9 recordings.

## RESULTS

The results for both approaches are displayed in Tables 2 and 3. As before, the accuracy represents the ratio of correctly classified windows to all windows. To account for the class imbalance, the mean recall and precision over all categories were added: the recall of a binary classifier is defined as the ratio of correctly classified positive to all true positive samples, whereas the precision of a binary classifier is defined as the ratio of correctly classified positive to all positively classified samples. These quantities for binary classifiers were extended to the multiclass problem by first calculating the respective quantity per class and then averaging the calculated quantities over all classes.

In both cases, the proposed FNNs outperformed the benchmark model. For the classic approach, the accuracy of the smaller FNN(3, 12) with 1,564 parameters had a 2.4%

Table 2. Comparison of the models’ quality under the classic approach.

Model	Accuracy	Recall	Precision
EEGNet	58.13	58.30	58.59
FNN(5, 10)	<b>61.60</b>	<b>61.81</b>	<b>61.49</b>
FNN(3, 12)	60.53	60.80	60.56

Table 3. Comparison of the models’ quality under the sliding windows approach.

Model	Accuracy	Recall	Precision
EEGNet	46.00	35.46	40.52
FNN(40, 20)	<b>50.69</b>	<b>42.47</b>	<b>45.50</b>
FNN(5, 10)	47.97	37.42	42.60

higher accuracy and the larger FNN(5, 10) with 2,344 trainable parameters had a 3.47% higher accuracy compared to the benchmark model with 2,526 parameters. The respective confusion matrices are displayed in Figure 8 of Appendix B.

For the sliding window approach, the FNN(5, 10) with a comparable number of parameters had a 1.97% higher accuracy compared to the benchmark model, whereas the larger FNN(40, 20) outperformed the benchmark model by 4.69% accuracy. Note that classes 4 - 7 are generally more difficult to detect. Yet, it can be seen from the confusion matrices in Figure 9 of Appendix B, that the classifications of the FNN(40, 20) are particularly better for those classes, which is also reflected in the recall and precision of the classifiers.

Both the FNNs and the (default) EEGNet are relatively simple models. It can be expected that the accuracies improve for both types of models if the hyperparameters are tuned carefully. Further improvements might be possible by changing the FNN’s architecture or simply using more layers.

## FULLY FUNCTIONAL PREDICTIONS

With the proposed methodology it is not only possible to predict scalar-valued labels and use the model for classification, but it is also possible to predict functional labels. With the sliding windows as before, we can try to predict the class label for each time point rather than one label for the whole window. This is particularly useful at the transition from one state to another because these transitions cannot be represented by a simple classification.

We trained a functional neural network with three functional convolutional layers as depicted in the bottom of Figure 1 to predict labels for each time instant of the sliding windows.

The overall accuracy was similar to the accuracy reported for the classifier above. In Figures 4 and 5 are the true and predicted labels of two windows at the transition from an inter-trial break (class 1) to the interval after a trial onset (class 3) and from the interval after a trial onset (class 3) to the interval after a cue onset (class 2). It can be seen from both figures that the predictions do not match the true labels perfectly and that the confidences at the border region are generally lower, but overall the predictions match the true labels.

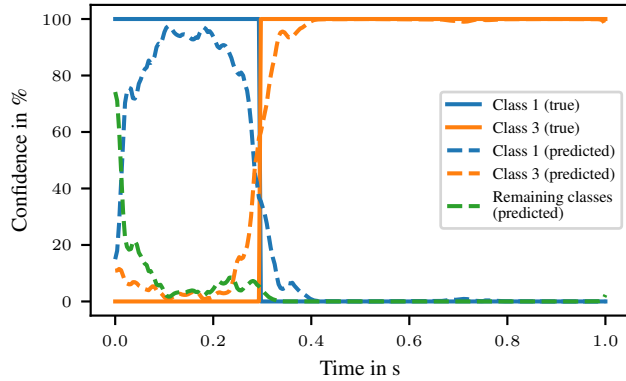


Figure 4. True and predicted class labels for a window of one second at the transition from class 1 (inter-trial break) to class 3 (time after trial onset).

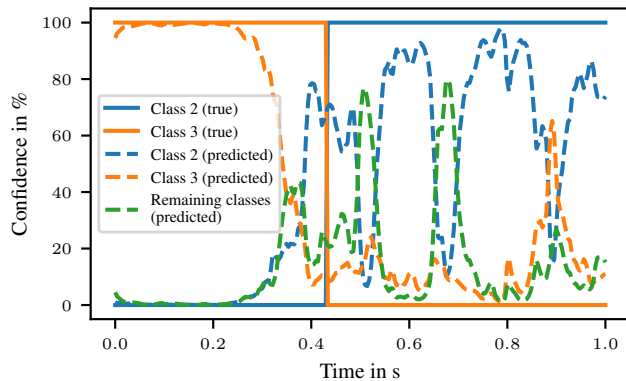


Figure 5. True and predicted class labels for a window of one second at the transition from class 3 (time after trial onset) to class 2 (time after cue onset).

## 6. Conclusions

In this work, we presented a new framework to analyze functional data with scalar- or functional-valued targets. We combined advantages of convolutional neural networks with those of functional data analysis. More specifically, we proposed a neural network that can be considered as shift

invariant while taking the intrinsic functional structure of the data into account.

We showed that even shallow models with only two convolutional and one dense layer are more powerful than the functional k-nearest neighbors algorithm for the simulated data. Further, the results of our case study suggest that FNNs with a similar amount of trainable weights outperform EEG-Net, the de facto standard model for the classification of EEG data.

The results of this paper suggest that functional neural networks are a relevant area for future research. First, it could be tested if FNNs work well with other types of time series and functional data, such as stock prices or temperature curves. It might be interesting to investigate if the methodology can be expanded to other types of data like images (considered as functions of the two variables height and width) or videos (considered as functions of the three variables time, height and width). Further, the choice of base is crucial for the performance of the network. Prior simulation studies showed that the Fourier base and Legendre polynomials lead to good results, but other bases might further improve the predictions. In FDA it is common to find roughness penalties as regularizers. Although a preliminary simulation study suggested that a base representation of the weight functions leads to better results, it would be interesting to study if the weight functions in the neural network can be learned directly while their smoothness would be ensured via corresponding roughness penalties. Finally, the proposed framework could be extended to other types of neural networks, such as recurrent neural networks or transformers.

## Acknowledgment

This work is supported by the Ministry of Economics, Innovation, Digitization and Energy of the State of North Rhine-Westphalia and the European Union, grant IT-2-2-023 (VAFES).

## References

- Azlan, W. A. W. and Low, Y. F. Feature extraction of electroencephalogram (EEG) signal-a review. In *2014 IEEE Conference on Biomedical Engineering and Sciences (IECBES)*, pp. 801–806. IEEE, 2014.
- Boonyakitanont, P., Lek-Uthai, A., Chomtho, K., and Songsiri, J. A review of feature extraction and performance evaluation in epileptic seizure detection using EEG. *Biomedical Signal Processing and Control*, 57: 101702, 2020.
- Boubchir, L., Daachi, B., and Pangracious, V. A review of feature extraction for EEG epileptic seizure detection



- and classification. In *2017 40th International Conference on Telecommunications and Signal Processing (TSP)*, pp. 456–460. IEEE, 2017.
- Bücher, A., Dette, H., and Heinrichs, F. A portmanteau-type test for detecting serial correlation in locally stationary functional time series. *arXiv preprint arXiv:2009.07312*, 2020.
- Cardot, H., Ferraty, F., and Sarda, P. Functional linear model. *Statistics & Probability Letters*, 45(1):11–22, 1999.
- Craik, A., He, Y., and Contreras-Vidal, J. L. Deep learning for electroencephalogram (EEG) classification tasks: a review. *Journal of neural engineering*, 16(3):031001, 2019.
- Cuevas, A., Febrero, M., and Fraiman, R. Linear functional regression: the case of fixed design and functional response. *Canadian Journal of Statistics*, 30(2):285–300, 2002.
- da Silva Lourenço, C., Tjepkema-Cloostermans, M. C., and van Putten, M. J. Efficient use of clinical EEG data for deep learning in epilepsy. *Clinical neurophysiology*, 132(6):1234–1240, 2021.
- Fan, J. and Gijbels, I. *Local polynomial modelling and its applications*. Routledge, 1 edition, 1996.
- Gabrys, R. and Kokoszka, P. Portmanteau test of independence for functional observations. *Journal of the American Statistical Association*, 102(480):1338–1348, 2007.
- Gasser, T. and Kneip, A. Searching for structure in curve samples. *Journal of the american statistical association*, 90(432):1179–1188, 1995.
- Kneip, A. and Gasser, T. Statistical tools to analyze data representing a sample of curves. *The Annals of Statistics*, pp. 1266–1305, 1992.
- Kokoszka, P. and Reimherr, M. *Introduction to functional data analysis*. Chapman and Hall/CRC, 2017.
- Lawhern, V. J., Solon, A. J., Waytowich, N. R., Gordon, S. M., Hung, C. P., and Lance, B. J. EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
- Olver, P. J. Introduction to the calculus of variations. *University of Minnesota*, 2016.
- Ramsay, J. O. and Li, X. Curve registration. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(2):351–363, 1998.
- Ramsay, J. O. and Silverman, B. W. *Functional Data Analysis*. Springer Science+Business Media, Inc., 2 edition, 2005.
- Rao, A. R. and Reimherr, M. Modern non-linear function-on-function regression. *arXiv preprint arXiv:2107.14151*, 2021a.
- Rao, A. R. and Reimherr, M. Non-linear functional modeling using neural networks. *arXiv preprint arXiv:2104.09371*, 2021b.
- Rossi, F. and Conan-Guez, B. Functional multi-layer perceptron: a non-linear tool for functional data analysis. *Neural networks*, 18(1):45–60, 2005.
- Rossi, F., Conan-Guez, B., and Fleuret, F. Functional data analysis with multi layer perceptrons. In *Proceedings of the 2002 International Joint Conference on Neural Networks. IJCNN’02 (Cat. No. 02CH37290)*, volume 3, pp. 2843–2848. IEEE, 2002.
- Rossi, F., Delannay, N., Conan-Guez, B., and Verleysen, M. Representation of functional data in neural networks. *Neurocomputing*, 64:183–210, 2005.
- Roy, Y., Banville, H., Albuquerque, I., Gramfort, A., Falk, T. H., and Faubert, J. Deep learning-based electroencephalography analysis: a systematic review. *Journal of neural engineering*, 16(5):051001, 2019.
- Sakoe, H. and Chiba, S. Dynamic programming algorithm optimization for spoken word recognition. *IEEE transactions on acoustics, speech, and signal processing*, 26(1): 43–49, 1978.
- Shang, H. L. A survey of functional principal component analysis. *AStA Advances in Statistical Analysis*, 98(2): 121–142, 2014.
- Silverman, B. W. *Density estimation for statistics and data analysis*. Routledge, 2018.
- Tangemann, M., Müller, K.-R., Aertsen, A., Birbaumer, N., Braun, C., Brunner, C., Leeb, R., Mehring, C., Miller, K. J., Mueller-Putz, G., et al. Review of the BCI competition IV. *Frontiers in neuroscience*, pp. 55, 2012.
- Wang, Q., Zheng, S., Farahat, A., Serita, S., and Gupta, C. Remaining useful life estimation using functional data analysis. In *2019 IEEE International Conference on Prognostics and Health Management (ICPHM)*, pp. 1–8. IEEE, 2019.
- Wang, Q., Wang, H., Gupta, C., Rao, A. R., and Khorasgani, H. A non-linear function-on-function model for regression with time series data. In *2020 IEEE International Conference on Big Data (Big Data)*, pp. 232–239. IEEE, 2020.

Zhao, J. Functional data learning by Hilbert feedforward neural networks. *Mathematical Methods in the Applied Sciences*, 35(17):2111–2121, 2012.

## A. Simulated Data - Samples

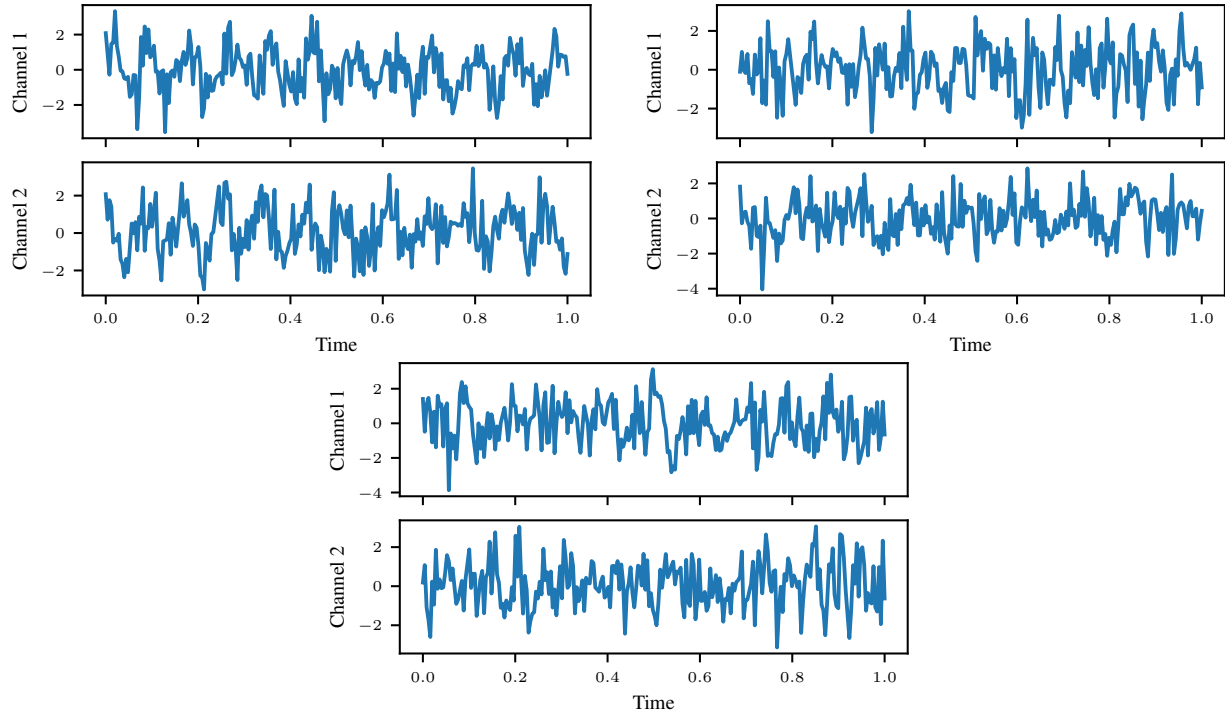


Figure 6. Samples from the simulated data set (I). Top left: Example of class 1 ( $\alpha$  frequencies only). Top right: Example of class 2 ( $\beta$  frequencies added in both channels, stronger in Channel 1). Bottom: Example of class 3 ( $\beta$  frequencies added in both channels, stronger in Channel 2).

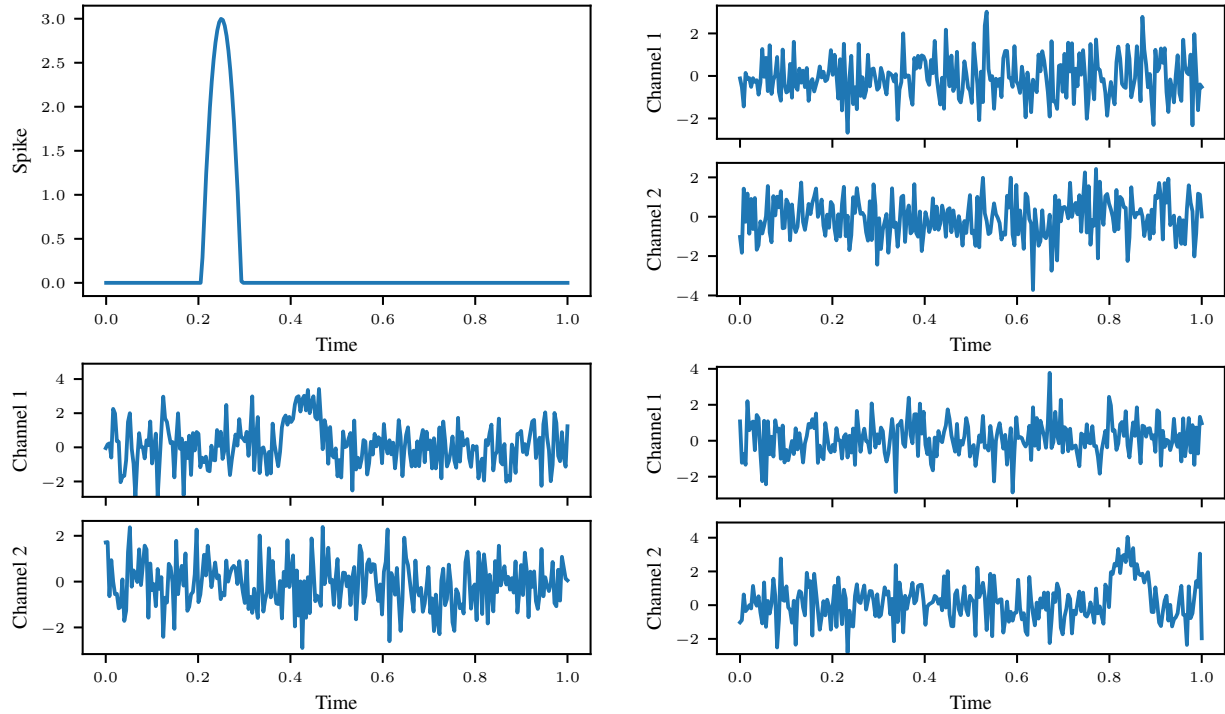


Figure 7. Samples from the simulated data set (II). Top left: “Spike” at  $t = 0.25$  with width  $w = 0.05$ . Top right: Example of class 1 (white noise in both channels). Bottom left: Example of class 2 (“spike” added to white noise in Channel 1). Bottom right: Example of class 3 (“spike” added to white noise in Channel 2).

## B. Real Data Experiments - Results

	1	2	3	4		1	2	3	4		1	2	3	4
1	13.8	4.7	3.8	1.5	1	13.6	2.7	4.6	2.9	1	13.7	2.8	3.9	3.3
2	3.0	16.0	3.8	2.2	2	2.2	18.0	2.4	2.2	2	2.6	17.2	2.6	2.6
3	5.4	4.7	13.6	3.2	3	3.9	3.4	13.8	6.0	3	4.3	3.2	12.9	6.6
4	3.0	3.7	2.9	14.8	4	2.8	2.2	3.0	16.2	4	2.8	1.7	3.1	16.7

Figure 8. Confusion matrices for different models under the classic approach. Left: EEGNet. Center: FNN(5, 10). Right: FNN(3, 12).

	1	2	3	4	5	6	7		1	2	3	4	5	6	7		1	2	3	4	5	6	7
1	20.8	0.5	3.6	0.3	0.4	0.2	0.4	1	19.9	0.3	3.0	0.7	0.8	0.6	0.7	1	20.9	0.3	3.5	0.3	0.4	0.2	0.5
2	1.5	6.0	3.3	0.4	0.7	0.2	0.5	2	0.9	8.0	1.8	0.4	0.5	0.4	0.4	2	1.3	7.4	2.5	0.3	0.3	0.2	0.3
3	8.0	1.6	13.2	0.4	0.6	0.3	0.7	3	6.2	1.0	13.6	0.9	1.2	0.7	1.0	3	7.8	1.3	13.5	0.5	0.5	0.4	0.8
4	3.1	0.6	1.8	1.3	1.5	0.4	0.8	4	3.0	0.2	1.4	2.3	1.2	0.8	0.7	4	3.6	0.4	1.7	1.4	1.0	0.6	0.8
5	2.5	0.6	1.8	0.6	1.9	0.2	0.5	5	2.6	0.3	1.1	0.8	2.3	0.6	0.6	5	3.2	0.3	1.8	0.5	1.5	0.4	0.5
6	3.0	0.9	2.8	0.5	0.6	1.0	1.2	6	2.6	0.4	1.6	0.9	1.1	2.0	1.4	6	3.4	0.6	2.2	0.6	0.6	1.1	1.3
7	2.9	0.5	2.6	0.4	0.3	0.6	1.8	7	2.4	0.3	1.4	0.5	0.8	1.1	2.7	7	3.4	0.4	2.0	0.3	0.4	0.7	2.1

Figure 9. Confusion matrices for different models under the sliding window approach. Left: EEGNet. Center: FNN(40, 20). Right: FNN(5, 10).