

---

# DECODING NEURAL SIGNALS WITH A COMPACT AND INTERPRETABLE CONVOLUTIONAL NEURAL NETWORK

---

A PREPRINT

Arthur Petrosyan, Mikhail Lebedev and Alexei Ossadtchi

June 2, 2020

## ABSTRACT

1 Brain-computer interfaces (BCIs) decode information from neural activity and send it to external  
2 devices. In recent years, we have seen an emergence of new algorithms for BCI decoding. Here  
3 we propose a compact architecture for adaptive decoding of electrocorticographic (ECoG) data into  
4 finger kinematics. We also describe a theoretically justified approach to interpreting the spatial  
5 and temporal weights in the architectures that combine adaptation in both space and time, such as  
6 ours. In these architectures the weights are optimized not only for decoding of target signals but  
7 also for tuning away from the interfering sources, in both the spatial and the frequency domains.  
8 When applied to a dataset taken from the repository of Berlin BCI IV competition, our architecture  
9 outperformed the competition winners without the need for feature selection. Moreover, by looking  
10 at the architecture weights we could explain in physiological terms how our algorithm decodes  
11 spatial and temporal parameters of finger kinematics. As such, the proposed architecture offers a  
12 good decoder and a tool for investigating neural mechanisms of motor control.

13 **Keywords:** *ECoG, limb kinematics decoding, deep learning, machine learning, weights interpretation, spatial filter,*  
14 *temporal filter*

## 15 1 Introduction

16 Brain-computer interfaces (BCIs) link the nervous system to external devices [4] or even the other brains [15]. While  
17 there exist many applications of BCIs [1], clinically relevant BCIs have received most attention that aid in rehabilitation  
18 of patients with sensory, motor, and cognitive disabilities [12]. Clinical uses of BCIs range from assistive devices to  
19 neural prostheses that restore functions abolished by neural trauma or disease [2].

20 BCIs can deal with a variety of neural signals [14, 8] such as, for example, electroencephalographic (EEG) potentials  
21 sampled with electrodes placed on the surface of the head [11], or neural activity recorded invasively with the elec-  
22 trodes implanted in the cortex [6] or places onto the cortical surface [18]. The latter method, which we consider here,  
23 is called electrocorticography (ECoG). Accurate decoding of neural signals is key to building efficient BCIs.

24 BCI signal processing comprises several steps, including signal conditioning, feature extraction, and decoding. In  
25 the modern machine-learning algorithms, feature extraction and decoding are not separate but rather simultaneous  
26 computations performed with the computational architectures called Deep Neural Networks (DNN) [9]. DNNs de-  
27 rive features automatically when executing regression or classification tasks. While it is often difficult to interpret  
28 the computations performed by a DNN, such interpretations are essential to gain understanding of the properties of  
29 brain activity contributing to decoding, and to ensure that artifacts do not affect the decoding results. In particular,  
30 interpretation of features computed by the first several layers of a DNN could shed light on the neurophysiological  
31 mechanisms underlying the behavior being studied. Ideally, by examining DNN weights, one should be able to match  
32 the algorithm's operation to the functions and properties of the neural circuitry to which the BCI connects. Moreover,  
33 we suggest that physiologically tractable DNN architectures could facilitate the development of efficient and versatile  
34 BCIs.

35 Several useful and compact architectures have been developed for processing EEG and ECoG data. The operation  
36 of some blocks of these architectures can be straightforwardly interpreted. Thus, EEGNet [7] contains explicitly  
37 delineated spatial and temporal convolutional blocks. This architecture yields high decoding accuracy with a minimal  
38 number of parameters. However, due to the cross-filter-map connectivity between any two layers, a straightforward  
39 interpretation of the weights is difficult. Some insight regarding the decision rule can be gained using DeepLIFT  
40 combined with the analysis of the hidden unit activation patterns. Schirrneister et al. describe two architectures:  
41 DeepConvNet and its compact version ShallowConvNet. The latter architecture consists of just two convolutional  
42 layers that perform temporal and spatial filtering, respectively [19].

43 Here we propose several novel approaches for making the operation of deep architectures tractable and interpretable  
44 neurophysiologically. Our approaches bear a resemblance with the recent study of Zubarev et al. [23] reporting two  
45 compact neural network architectures, LF-CNN and VAR-CNN, that outperformed the other decoders of MEG data,  
46 including linear models and more complex neural networks such as ShallowFBCSP-CNN, EEGNet-8 and VGG19. LF-  
47 CNN and VAR-CNN contain only a single non-linearity, which distinguishes them from most other DNNs. Because  
48 of this feature, the weights of these architectures are readily interpretable with the well-established approaches for  
49 interpreting the weights in linear models. Specifically, the spatial weights can be interpreted based on the principles of  
50 the estimation theory [5] combined with several additional assumptions, like the network training provides a Wiener-  
51 optimal solution and the subsequent temporal filtering can be disregarded. As to the temporal convolution weights,  
52 they can be interpreted by considering the Fourier-domain representations (with a caveat that the input data spectral  
53 characteristics are not taken into account).

54 While the compact architecture described here is conceptually similar to LF-CNN, our goals were different from  
55 those of Zubarev et al. [23]. We developed a theoretically justified approach to interpreting the spatial and temporal  
56 weights. Our method applies the optimal estimation theory to the space of temporally embedded multichannel data,  
57 with factorized spatial and temporal processing. This method allows us to consider two factors: (1) the weights  
58 optimizing the output correlation with the target signal, and (2) the weights minimizing the interference from the  
59 sources in both spatial and frequency domains.

## 60 2 Methods

61 Figure 1 illustrates the relationship between motor behavior (hand movements), brain activity, and ECoG recordings.  
 62 The activity,  $e(t)$ , of a set of neuronal populations,  $G_1 - G_I$ , engaged in motor control, is converted into a movement  
 63 trajectory,  $z(t)$ , through a non-linear transform  $H$ :  $z(t) = H(e(t))$ . The activity of populations  $A_1 - A_J$  is unrelated to  
 64 movement. The recordings of  $e(t)$  with a set of sensors are represented by a  $K$ -dimensional vector of sensor signals,  
 65  $\mathbf{x}(t)$ . This vector can be modeled as a linear mixture of signals resulting from the application of forward-model  
 66 matrices  $\mathbf{G}$  and  $\mathbf{A}$  to task-related sources,  $\mathbf{s}(t)$ , and task-unrelated sources,  $\mathbf{f}(t)$ , respectively:

$$\mathbf{x}(t) = \mathbf{G}\mathbf{s}(t) + \mathbf{A}\mathbf{f}(t) = \sum_{i=1}^I \mathbf{g}_i s_i(t) + \sum_{j=1}^J \mathbf{a}_j f_j(t) \quad (1)$$

67 We will refer to the noisy component of the recording as  $\eta(t) = \sum_{j=1}^J \mathbf{a}_j f_j(t)$ .

68 Linear inverse mapping is commonly used to derive the activity of sources from the sensor signals:  $\hat{\mathbf{s}}(t) = \mathbf{W}^T \mathbf{X}(t)$ ,  
 69 where columns of  $\mathbf{W}$  form a spatial filter that counteracts the volume conduction effect and decreases the effect of  
 70 noisy sources.

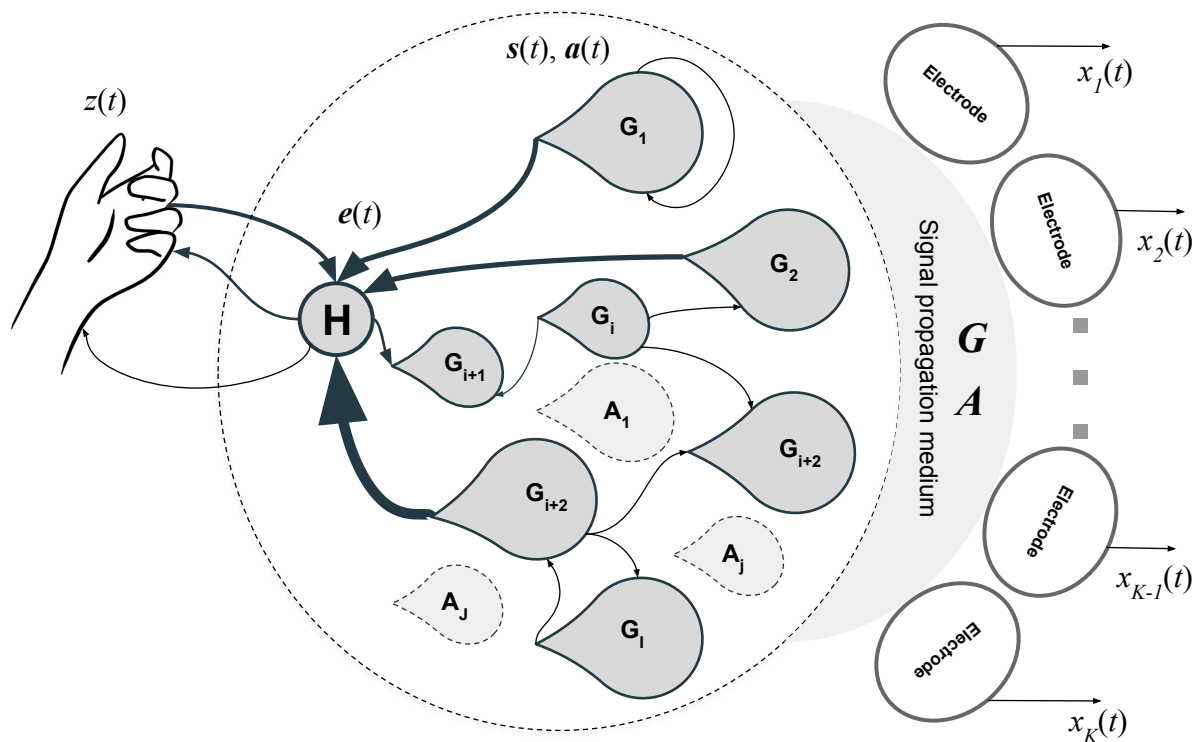


Figure 1: Phenomenological diagram

71 Neuronal correlates of motor planning and execution have been extensively studied [22]. In the cortical-rhythm do-  
 72 main, alpha and beta components of the sensorimotor rhythm envelope desynchronize just prior to the execution of a  
 73 movement and rebound with a significant overshoot upon the completion of a motor act [13]. The magnitude of these

74 modulations correlates with the person’s ability to control a motor-imagery BCI [16]. Additionally, the frequency of  
 75 beta bursts in the primary somatosensory cortex is inversely correlated with the ability to detect tactile stimuli [20].  
 76 Intracranial recordings, such as ECoG, allow reliable measurement of the faster gamma band activity, which is tempo-  
 77 rally and spatially specific to movement patterns [21]. Overall, rhythmic components of brain sources,  $s(t)$ , appear to  
 78 be useful for BCI implementations. These rhythmic signals can be computed as linear combinations of band-passed  
 79 sensor data,  $\mathbf{x}(t)$ .

80 The most straightforward approach for extracting the kinematics,  $z(t)$ , from brain recordings,  $\mathbf{x}(t)$ , is to directly  
 81 learn the mapping  $z(t) = \mathcal{F}(\mathbf{x}(t))$ . To do so, one needs to parametrically describe this mapping. Here we used a  
 82 specific network architecture for this purpose. The architecture was constructed in a close correspondence with the  
 83 neurophysiological description of the observed phenomena, which facilitated our ability to interpret the results.

## 84 2.1 Network architecture

85 The compact and adaptable architecture that we developed is shown in Figure 2. An adaptive envelope  
 86 extractor is the key component of this architecture. The envelope extractor, which is a module widely used in signal  
 87 processing systems, was implemented using modern DNN primitives, namely a pair of convolutional operations that  
 88 perform band-pass and low-pass filtering and one non-linearity ReLu(-1) that corresponds to computing the absolute  
 89 value of the output of the first 1-D convolutional layer. To make the decision rule of this structure tractable, we used  
 90 non-trainable batch normalization when streaming the data through the structure. All input signals were standardized.

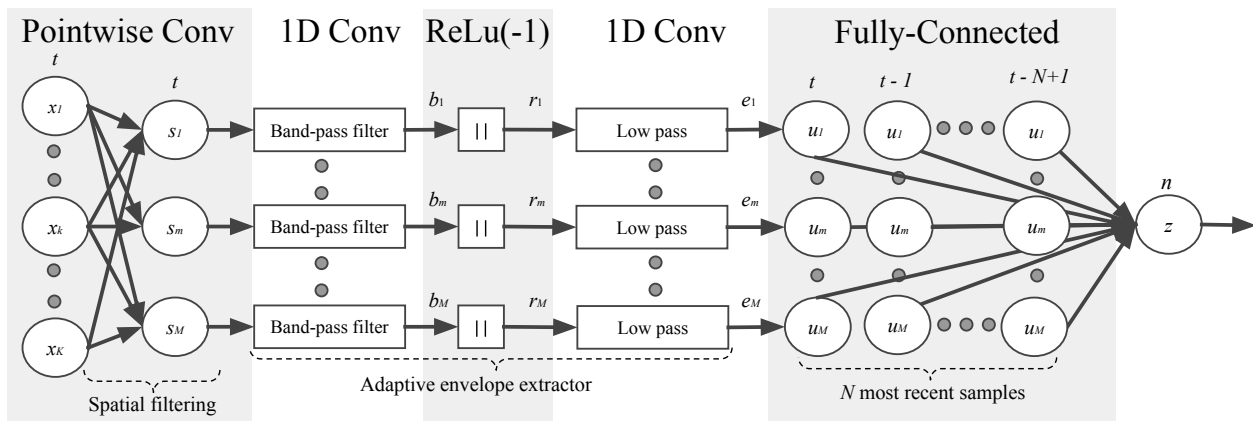


Figure 2: The proposed compact DNN architecture

91 In our architecture, the envelope detectors received spatially filtered sensor signals  $s_m$  that were calculated by the  
 92 point-wise convolutional layer. This layer counteracted the volume-conduction processes represented by the forward-  
 93 model matrix  $\mathbf{G}$  in our phenomenological model (Figure 1). Next, we approximated operator  $H$  as a function of the  
 94 lagged instantaneous power of the narrow-band source timeseries. This was done with a fully connected layer that  
 95 mixed the samples of envelopes,  $e_m(n)$ , into a single estimate of the kinematic parameter,  $z(n)$ .

## 96 2.2 Two regression problems and DNN weights interpretation

97 The proposed architecture processes data in chunks of a prespecified length of  $N$  samples. The processing of a chunk  
 98 of input data,  $\mathbf{X}(t) = [\mathbf{x}(t), \mathbf{x}(t-1), \dots, \mathbf{x}(t-N+1)]$ , by the first two layers performing spatial and temporal filtering  
 99 can be described for the  $m$ -th element as

$$b_m(n) = \mathbf{w}_m^T \mathbf{X}(t) \mathbf{h}_m \quad (2)$$

100 The non-linearity,  $ReLU(-1)$ , in combination with the low-pass filtering performed by the second convolutional layer,  
 101 extracts the envelopes of rhythmic signals.

102 The analytic signal is mapped one-to-one to its envelope [3]. Additionally, for the original real-valued data, the  
 103 imaginary part of the analytic signal is uniquely determined as Hilbert transform. Therefore, the adjustment of the  
 104 spatial and temporal filter weights to obtain some specific envelope  $e_m(t)$  is equivalent to the adjustment of the weights  
 105 to obtain this envelope's generating analytic signal  $s_m(t)$ . Accordingly, this is a linear regression problem where either  
 106 spatial or temporal weights are fixed and temporal or spatial weights are sought, correspondingly.

107 We assume that training of the adaptive envelope detectors results in optimal spatial and temporal convolution weights,  
 108  $\mathbf{w}_m^*$  and  $\mathbf{h}_m^*$ , correspondingly. Then, optimal spatial filter weights can be obtained as a solution to a convex optimiza-  
 109 tion problem formulated over spatial or temporal subset of parameters:

$$\mathbf{w}_m^* = \underset{\mathbf{w}_m}{\operatorname{argmin}} \{ \| b_m(n) - \mathbf{w}_m^T \mathbf{X}(t) \mathbf{h}_m^* \|_2^2 \} \quad (3)$$

110 where the temporal weights are fixed at their optimal value,  $\mathbf{h}_m^*$ . Similarly, when spatial weights are fixed at the  
 111 optimal value  $\mathbf{w}_m^*$ , temporal weights are expressed by the equation:

$$\mathbf{h}_m^* = \underset{\mathbf{h}_m}{\operatorname{argmin}} \{ \| b_m(t) - \mathbf{w}_m^{*T} \mathbf{X}(t) \mathbf{h}_m \|_2^2 \} \quad (4)$$

112 Given the forward model 1 and the regression problem 3 and assuming statistical independence of the rhythmic poten-  
 113 tials  $s_m(t)$ ,  $m = 1, \dots, M$ , the topographies of the underlying neuronal populations can be found as [? 5]

$$\mathbf{g}_m = E\{\mathbf{Y}(t)\mathbf{Y}^T(t)\} \mathbf{w}_m^* = \mathbf{R}_m^Y \mathbf{w}_m^* \quad (5)$$

114 where  $\mathbf{Y}(t) = \mathbf{X}(t)\mathbf{h}_m$  is a temporally filtered chunk of multichannel data and  $\mathbf{R}_m^Y = E\{\mathbf{Y}(t)\mathbf{Y}^T(t)\}$  is a  $K \times K$   
 115 covariance matrix of the temporally filtered data, assuming that  $x_k(t)$ ,  $k = 1, \dots, K$  are all zero-mean processes.  
 116 Thus, when interpreting individual spatial weights corresponding to each of the  $M$  paths of the architecture shown  
 117 in Figure 2 one has to take into account the temporal filter weights  $\mathbf{h}_m$  to which the individual  $m$ -th branch is  
 118 tuned. Therefore, to transform the spatial weights of different branches into spatial patterns, branch-specific covariance  
 119 matrices  $\mathbf{R}_m^Y$  should be used that depend on the temporal convolution weights of each particular branch.

120 The temporal weights can be interpreted in a similar way. The temporal pattern is calculated as

$$\mathbf{q}_m = E\{\mathbf{V}(t)\mathbf{V}^T(t)\}\mathbf{h}_m^* = \mathbf{R}_m^V\mathbf{h}_m^* \quad (6)$$

121 where  $\mathbf{V}(t) = \mathbf{X}^T(t)\mathbf{w}_m^*$  is a spatially filtered chunk of incoming data and  $\mathbf{R}_m^V = E\{\mathbf{V}(t)\mathbf{V}^T(t)\}$  is an  $N \times N$  covari-  
 122 ance matrix of the spatially filtered data, assuming that  $x_k(t)$ ,  $k = 1, \dots, K$  are all zero-mean processes. As with the  
 123 spatial patterns, when interpreting individual temporal weights corresponding to each of the  $M$  branches of the archi-  
 124 tecture shown in Figure 2, one has to take into account the spatial filter weights  $\mathbf{h}_m$  used to filter the individual  $m$ -th  
 125 branch. To transform the temporal convolution weights of different branches into temporal patterns, branch-specific  
 126 covariance matrices  $\mathbf{R}_m^V$  should be used that depend on the spatial convolution weights of each particular branch. To  
 127 assess the temporal pattern, we usually explore it in the frequency domain, i.e.  $Q_m(f) = \sum_{t=0}^{t=N-1} q_m(t)e^{-j2\pi ft}$ ,  
 128 where  $q_m(t)$  is the  $t$ -th element of temporal pattern vector  $\mathbf{q}_m$ .

129 Hitherto, we assumed that data chunk length  $N$  is equal to the length of the filters in the first convolutional layer. In  
 130 general, this does not have to be the case. Our assumption emphasizes the formal similarity between the spatial and  
 131 temporal dimensions. Additionally, we emphasize that the approach to the interpretation of temporal patterns requires  
 132 taking into account the correlation structure of the independent variable in the regression model. When the data chunk  
 133 is longer than the filter length, equation 2 has to be rewritten using the convolution operation. In this case, instead of  
 134 a scalar, the equation returns a vector of samples, with the vector length depending on the choice of strategy used to  
 135 deal with the transient at the edges of the chunk. It is also easier to operate in the frequency domain from the very  
 136 beginning, and use the standard Wiener filtering arguments. In the frequency domain, the Wiener filter weights can be  
 137 expressed as a function of the power spectral density  $P_m^{yy}(f)$  of the spatially filtered sensor data  $y_m(t) = \mathbf{w}_m^T \mathbf{x}(t)$  in  
 138 the  $m$ -th branch and the density of cross-spectrum,  $P_m^{sy}(f)$ , between  $s_m(t)$  and  $y_m(t)$  :

$$Q_m^*(f) = \frac{P_m^{sy}(f)}{P_m^{yy}(f)} \quad (7)$$

139 Then, using the assumption that  $\eta(t)$  and  $\mathbf{s}(t)$  in 1 are statistically independent, we obtain the following expressions:

$$H_m^*(f) = \frac{P_m^{ss}(f)}{P_m^{ss}(f) + P_m^{\eta\eta}(f)} = \frac{P_m^{ss}(f)}{P_m^{yy}(f)} \quad (8)$$

140 Therefore, the frequency-domain pattern of the signal isolated by the  $m$ -th branch spatial filter can be computed as

$$Q_m^*(f) = P_m^{yy}(f)H_m^*(f) \quad (9)$$

141 where  $H_m^*(f)$  in 9 is the Fourier transform of the vector  $\mathbf{h}_m^*$  containing temporal-convolution weights identified during  
 142 the adaptation of the envelope detector in the  $m$ -th branch. Viewing this result as a product of learning, it means that  
 143 the learned vector of temporal convolution weights,  $\mathbf{h}_m$ , represents the power spectral density of the brain potentials  
 144 that are important for decoding sensor signals,  $\mathbf{x}(t)$ , into the kinematics,  $z(t)$ .

145 The spatial patterns of neuronal sources recovered from the spatial filtering weights are routinely used for dipole fitting  
 146 to localize functionally important neural sources. The temporal patterns interpreted according to 9 and 6 can be used  
 147 to fit the models of neural population dynamics, which are relevant to specific decoding tasks.

## 148 2.3 Simulations

149 To tackle the performance of the proposed architecture, we performed a set of simulations. The simulated data corre-  
150 sponded to the setting shown in the phenomenological diagram (Figure 1). We simulated  $I = 4$  task related sources  
151 with rhythmic potentials,  $s_i(t)$ . The potentials of these four task related populations were generated as narrow-band  
152 processes (in 30-80 Hz], 80-120 Hz], 120-170 Hz and 170-220 Hz bands) resulting from filtering of Gaussian pseudo-  
153 random sequences with a bank of FIR filters. We then simulated the kinematics,  $z(t)$ , as a linear combination of  
154 the four envelopes of these rhythmic signals with randomly generated vector of coefficients. We used  $J = 40$  task-  
155 unrelated rhythmic sources with activation timeseries obtained similarly to the task-related sources but with filtering  
156 within the following four bands: 40-70 Hz, 90-110 Hz, 130-160 Hz, and 180 - 210 Hz bands. As a result, we obtained  
157 10 task unrelated sources active in each of these bands making it a total of  $J = 40$  task unrelated sources. To simulate  
158 volume conduction effect, we randomly generated a  $4 \times 5$  dimensional forward matrix  $\mathbf{G}$  and a  $40 \times 5$  dimensional  
159 forward matrix  $\mathbf{A}$ . These matrices mapped the task-related and task-unrelated activity, respectively, onto the sensor  
160 space.

161 We generated 15 minutes worth of data sampled at 1,000 Hz and split them into two equal contiguous parts. We used  
162 the first part for training and the second for testing

## 163 3 Experimental datasets

164 To compare the performance of our simple neural network with the top linear models that rely on preset features,  
165 we used publicly available data collected by Kubanek et al [ ] from the BCI Competition IV. This dataset contains  
166 concurrent multichannel ECoG and finger flexion measurements collected in three epileptic patients implanted with  
167 ECoG electrodes for medical reasons. The database consists of 400 s of training data and 200 s of test data. The  
168 recordings were conducted with 64 or 48 electrodes placed over the sensorimotor cortex. The exact spatial locations  
169 and the order of the electrodes are not provided. As a baseline in this comparison, we chose the winning solution  
170 offered by Nanying Liang and Laurent Bougrain [10]. This solution employs extracting the amplitudes of the data  
171 filtered in 1-60 Hz, 60-100 Hz, and 100-200 Hz band followed by a pairwise feature selection and decoded using  
172 Wiener filter with  $N = 25$  taps from the immediate past.

173 The other dataset comes from our laboratory. The recordings were conducted with a 64-channel Adtech microgrid  
174 connected to EB Neuro BE Plus LTM Bioelectric Signals Amplifier System that sampled data at 2048 Hz. The am-  
175 plifier software streamed data via Lab Streaming Layer protocol. The experimental software supported this protocol,  
176 implemented the experimental paradigm (a finger movement task) and synced ECoG and kinematics. Finger kinemat-  
177 ics was captured by Perception Neuron system as relative angles for the sensor units attached to finger phalanges, and  
178 sampled at 120 Hz. Finger flexion-extension angle was used as kinematics timeseries,  $z(t)$ .

179 Recordings were obtained in two patients with pharmaco-resistant form of epilepsy; ECoG electrodes were implanted  
180 for the purpose of pre-surgery localization of epileptic foci and mapping of eloquent cortex. Thus, for these data, unlike  
181 Kubanek [ ], we knew cortical location of each electrode and could visualize spatial patterns of activity with high  
182 accuracy. The patients performed self-paced flexions of each individual finger for 1 min. The study was conducted  
183 according to the ethical standards of the 1964 Declaration of Helsinki. All participants provided written informed

184 consent prior to the experiments. The ethics research committee of the National Research University, The Higher  
185 School of Economics approved the experimental protocol of this study

## 186 4 Results for simulated data

### 187 4.1 Adaptive envelop detector

188 As described in Methods, to interpret optimal temporal convolution weight,s we need to consider the spectral character-  
189 istics of neural recordings. To illustrate this, we trained a single-channel adaptive envelope detector in the environment  
190 with the interference occupying a subrange of the target signal band. As can be seen from Figure 3, the Fourier profile  
191 of the identified temporal convolution weights can not be used to assess the power spectral density of the underlying  
192 signal. At the same time, the expression in 8 allows us to obtain a proper pattern that matches well the simulated spec-  
193 tral profile. Conversely, using the FFT of the convolutional filter weights yields fundamentally erroneous estimates of  
194 the frequency-domain patterns and erroneous interpretation of the underlying neurophysiology.

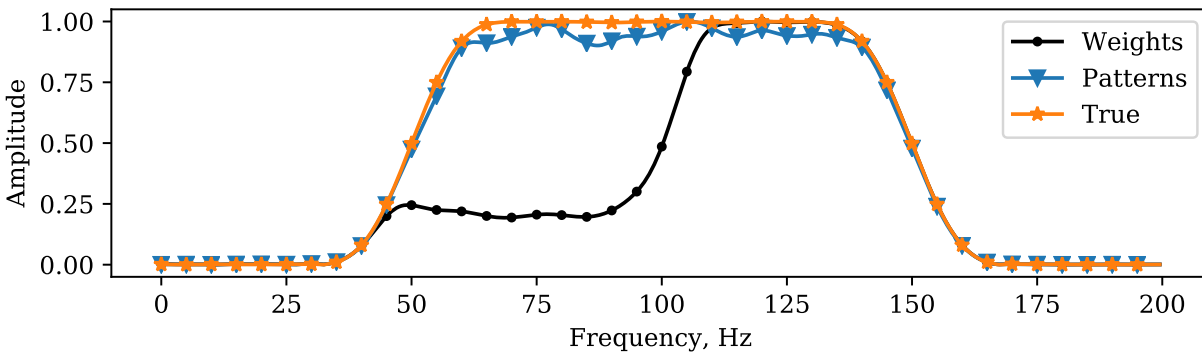


Figure 3: The importance of taking into account the power spectral density of the independent variable when interpreting linear regression weights. The true pattern (\*) gets erroneously reproduced as the FFT of the temporal convolution weights (●). Taking into account the power spectral density of the spatially filtered signal allows to fix this situation (▼)

### 195 4.2 Realistic simulations

196 For the simulated data, we trained the algorithm to predict the kinematic variable  $z(t)$ . In the noiseless case, the  
197 proposed architecture achieved accuracy of 99% measured as correlation coefficient between the true and recovered  
198 kinematics (Figure 4). We then compared the envelopes at each of the four branches of our architecture and observed  
199 that the true latent variable timeseries (in the form of the underlying narrow-band envelopes) matched very well those  
200 estimated with our architecture (Figure 5). The correlation between the estimated and true envelope timeseries fell  
201 into the 87 - 96 % range.

202 As described in Methods, for spatial weights interpretation, we used the linear estimation theoretic approach [5]. To  
203 warn against its naive implementation in the context of architectures that combine spatial and temporal filtering, we  
204 computed spatial patterns where we used the input data covariance,  $\mathbf{R}^X$ , without taking into account the individual-  
205 branch temporal filters. In the corresponding plots, we refer to the patterns determined using this approach as *Patterns*



206 *vanilla*. The proper way to apply this estimation approach is to compute spatial covariance,  $\mathbf{R}^Y$ , for the temporally  
207 filtered data  $\mathbf{6}$ . These properly determined patterns are labeled as *Patterns vanilla*.

208 In the right column of Figure 6, we show the results for the noiseless case for all four branches of the network. As  
209 expected, the spatial *Patterns vanilla* and *Patterns* are identical and match the ground truth exactly. The left column  
210 shows Fourier representations of the temporal weights where we can observe that in the noise-free scenario Fourier  
211 representations of the temporal weights matches exactly the power spectral density of the simulated data.

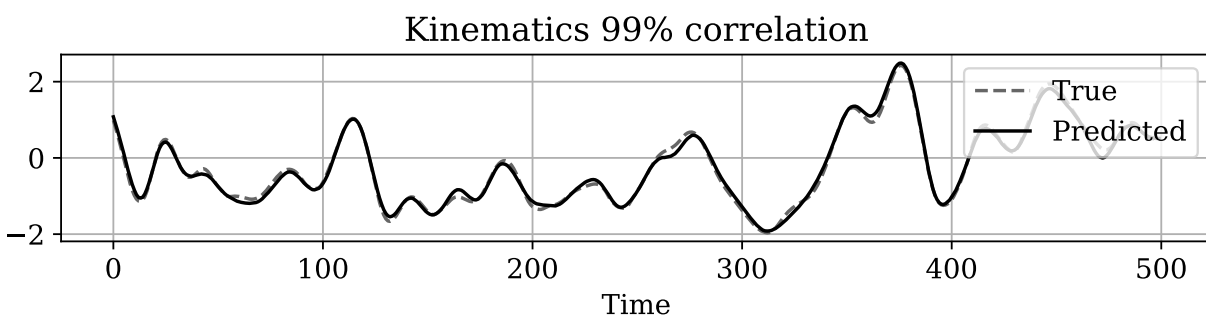


Figure 4: Realistic simulations. Actual and decoded kinematics.

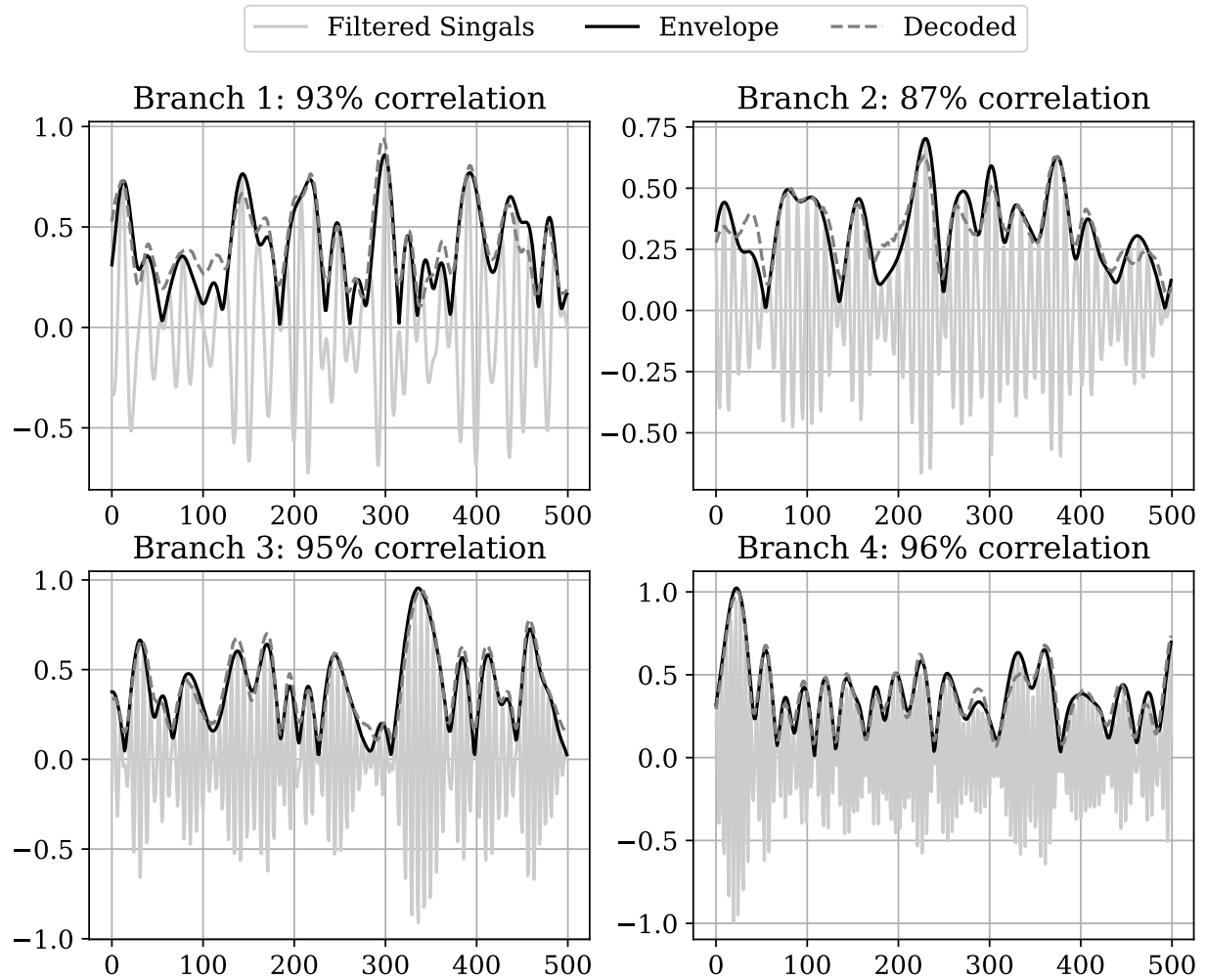


Figure 5: Branch envelopes decoding. Comparison between the true and the decoded envelope.

212 In the noisy case demonstrated in Figure 7, only *Patterns vanilla* match well with the simulated topographies of the  
213 underlying sources. Spectral characteristics of the trained temporal filtering weights exhibit characteristic dips in the  
214 bands corresponding to the activity of the interfering sources. After applying the theoretic estimation 9, we obtain the  
215 spectral patterns that more closely match the simulated ones and have the dips compensated.

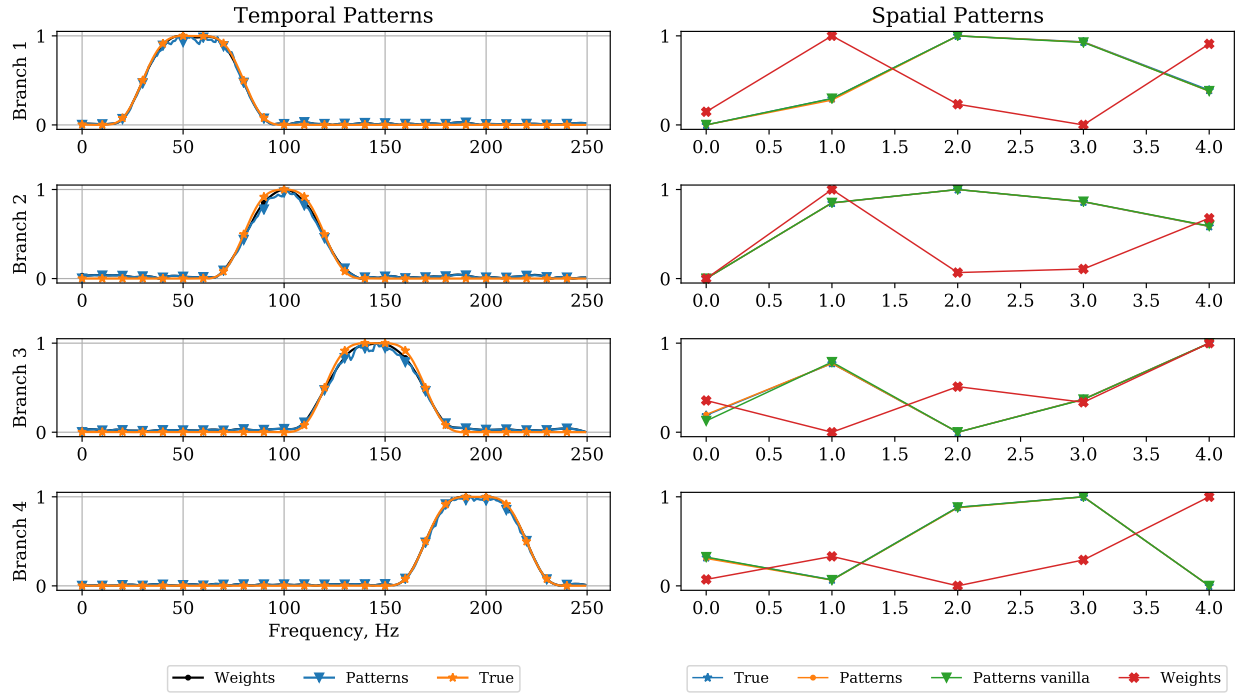


Figure 6: Temporal (left) and spatial(right) patterns obtained for the noiseless case. See the main text for description.

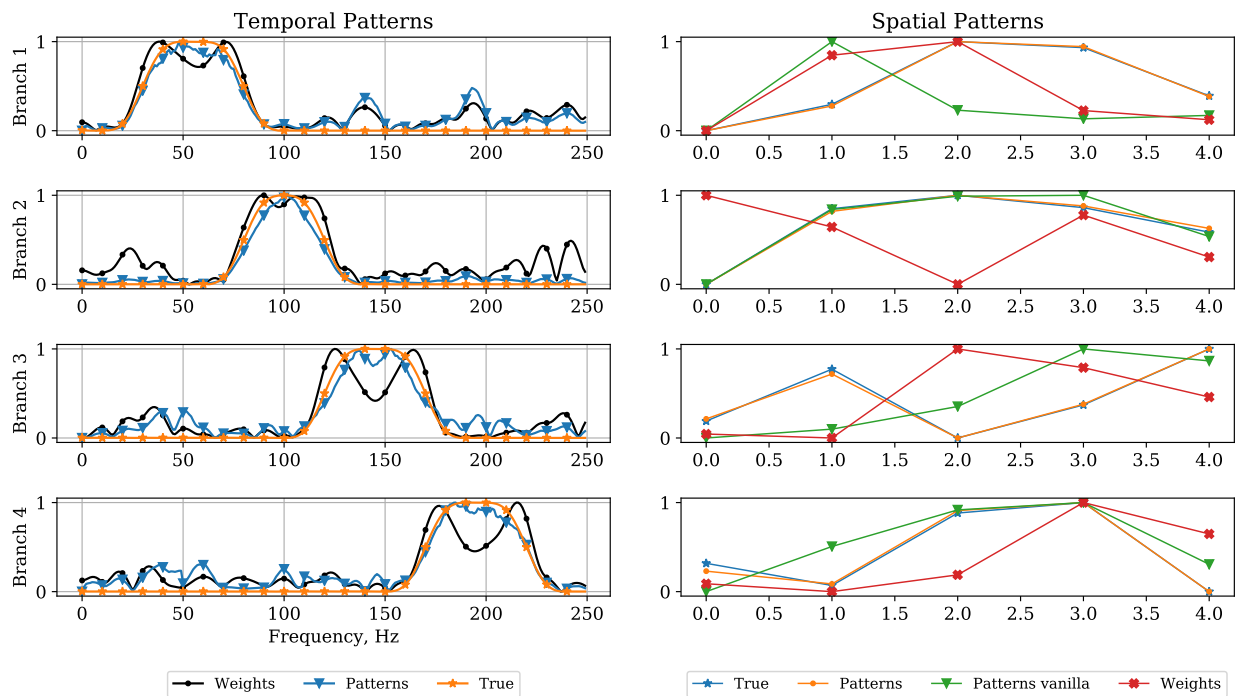


Figure 7: Temporal (left) and spatial(right) patterns obtained for the noisy case, SNR = 1.5. See the main text for description.

## 216 5 Analysis of experimental data

### 217 5.1 Berlin BCI Competition IV data

218 In the context of electrophysiological data processing, the major advantage of the architectures inspired by the deep-  
219 learning principle is their ability to automatically select features while performing classification or regression tasks  
220 [17]. When applied to the data from Berlin BCI Competition IV, our architecture – based on the adaptive envelope  
221 detectors – performed on par or better than the winning solution by Lian and Bougrain [10], see Table 1.

Subject 1					
	Thumb	Index	Middle	Ring	Little
Winner	0,58	0,71	0,14	0,53	0,29
NET	0,54	0,7	0,2	0,58	0,25

Subject 2					
	Thumb	Index	Middle	Ring	Little
Winner	0,51	0,37	0,24	0,47	0,35
NET	0,5	0,36	0,22	0,4	0,23

Subject 3					
	Thumb	Index	Middle	Ring	Little
Winner	0,59	0,51	0,32	0,53	0,2
NET	0,71	0,48	0,5	0,52	0,61

Table 1: Comparison of the performance of the proposed architecture (NET) and the winning solution (Winner) of Berlin BCI IV competition dataset(4)

### 222 5.2 The CBI data

223 We also applied the proposed solutions to the recordings that we conducted in two patients implanted with  $8 \times 8$  ECoG  
224 grids over the sensorimotor cortex.

225 The following table shows the accuracy achieved with the proposed architecture for the decoding of finger movements.

	Thumb	Index	Ring	Little
Subject 1	0,48	<b>0,79</b>	0,61	0,32
Subject 2	0,73	0,55	0,78	<b>0,79</b>

Table 2: Decoding performance achieved in the two CBI patients. The table show correlation coefficient between the actual and the decoded finger trajectory for the four fingers in two patients.

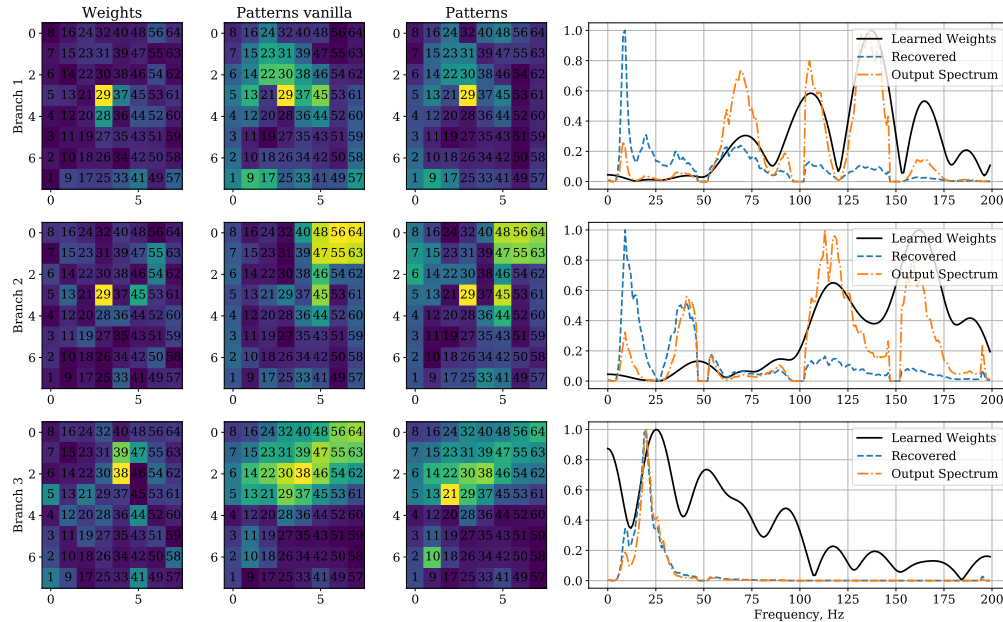


Figure 8: Network weights interpretation for the index finger decoder in CBI patient 1. Each row of plots corresponds to one of the three branches of the trained decoder. The left most column shows color-coded spatial filter weights, next two columns correspond to vanilla and properly recovered spatial patterns. The fourth column interprets temporal filter weights in the Fourier domain. Filter weights - solid line, power spectral density (PSD) pattern of the underlying LFP - dashed line. Another dashed line more similar to the filter weights Fourier coefficients is the PSD of the signal at the output of the temporal convolution block.

226 Figures 8 and 9 depict the interpretation of the obtained spatial and temporal weights. The plots are shown for the  
 227 finger with the highest decoding accuracy (highlighted in bold in Table 2) for two patients.

228 The decoding architecture for both patients had three branches and each branch was tuned to a source with specific  
 229 spatial and temporal patterns. In Figure 8, we show the spatial filter weights, vanilla patterns and proper patterns  
 230 interpreted using the expression described in the the Methods section. It can be seen that, while the temporal filter  
 231 weights (solid line) clearly emphasized the frequency range above 100 Hz in the first two branches, the actual spectral  
 232 pattern of the source (dashed line) in addition to the gamma-band content had a peak at around 11 Hz (1st, 2nd  
 233 branches) and in the 25-50 Hz range (2nd branch). These peaks likely correspond to the sensorimotor rhythm and  
 234 low-frequency gamma rhythms, respectively.

235 The third branch appears to capture the lower-frequency range and its spatial pattern is noticeably more diffuse than  
 236 that in the first two branches that capture the higher-frequency components. This is consistent with the phenomenon  
 237 that the size and activation frequency of neuronal populations are reciprocally proportional.

238 Similar observations can be made from Figure 9 that shows to the decoding results for the little finger in patient 2.

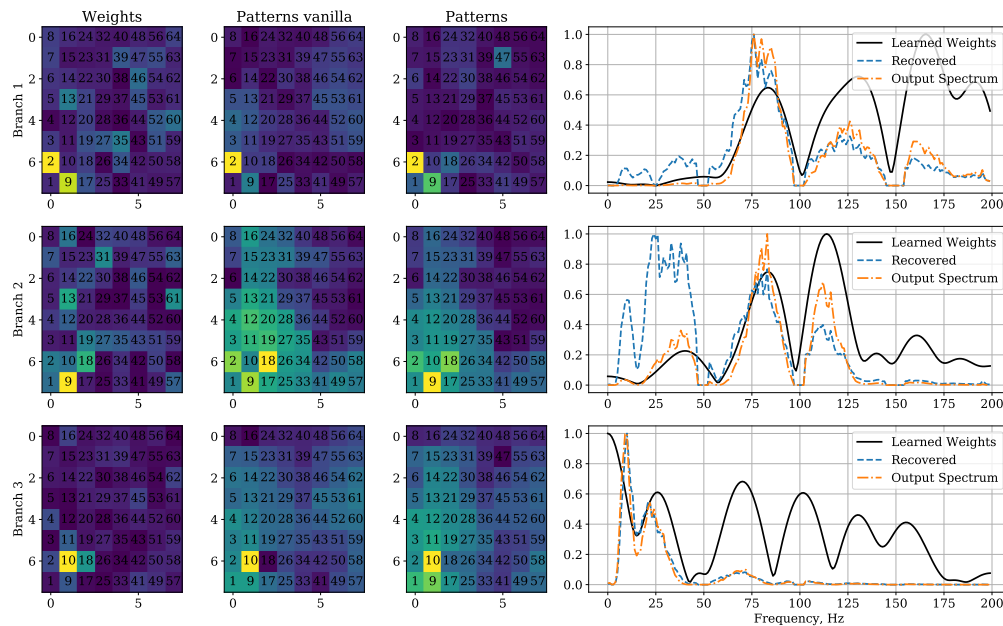


Figure 9: Network weights interpretation for the little finger decoder in CBI patient 2. Each row of plots corresponds to one of the three branches of the trained decoder. The left most column shows color-coded spatial filter weights, next two columns correspond to vanilla and properly recovered spatial patterns. The fourth column interprets temporal filter weights in the Fourier domain. Filter weights - solid line, power spectral density (PSD) pattern of the underlying LFP - dashed line. Another dashed line more similar to the filter weights Fourier coefficients is the PSD of the signal at the output of the temporal convolution block.

## 239 6 Conclusion

240 We developed a novel compact and neurophysiologically interpretable architecture. Using this architecture, we ex-  
 241 tended the weights interpretation approach previously applied in [5] to the interpretation of the temporal convolution  
 242 weights. We tested the proposed approach using simulated and experimental data. In the realistically simulated data,  
 243 our architecture recovered with high accuracy the neuronal substrate that contributed to the kinematics data.

244 We also applied the proposed architecture to an experimental dataset taken from the repository of Berlin BCI IV  
 245 competition. Our architecture delivered similar or better decoding accuracy as compared the winning solution of the  
 246 BCI competition [10]. In contrast to the traditional approaches, our architecture did not require any preset features.  
 247 Instead, after the architecture was trained to decode finger kinematics, we could interpret the weights and extracted  
 248 physiologically meaningful patterns corresponding to both spatial and temporal convolution weights.

## 249 7 Acknowledgment

250 This work is supported by the Center for Bioelectric Interfaces NRU HSE, RF Government grant, ag.  
 251 No.14.641.31.0003.

## 252 References

- 253 [1] Sarah N Abdulkader, Ayman Atia, and Mostafa-Sami M Mostafa. Brain computer interfacing: Applications and  
254 challenges. *Egyptian Informatics Journal*, 16(2):213–230, 2015.
- 255 [2] Ujwal Chaudhary, Niels Birbaumer, and Ander Ramos-Murguialday. Brain–computer interfaces for communica-  
256 tion and rehabilitation. *Nature Reviews Neurology*, 12(9):513, 2016.
- 257 [3] Stefan L Hahn. On the uniqueness of the definition of the amplitude and phase of the analytic signal. *Signal*  
258 *Processing*, 83(8):1815–1820, 2003.
- 259 [4] Nicholas G Hatsopoulos and John P Donoghue. The science of neural interface systems. *Annual review of*  
260 *neuroscience*, 32:249–266, 2009.
- 261 [5] Stefan Haufe, Frank Meinecke, Kai Gorgen, Sven Dahne, John-Dylan Haynes, Benjamin Blankertz, and Felix  
262 Biemann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage*,  
263 87:96–110, 2014.
- 264 [6] Mark L Homer, Arto V Nurmikko, John P Donoghue, and Leigh R Hochberg. Sensors and decoding for intracor-  
265 tical brain computer interfaces. *Annual review of biomedical engineering*, 15:383–405, 2013.
- 266 [7] Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J  
267 Lance. Eegnet: A compact convolutional network for eeg-based brain-computer interfaces. *arXiv preprint*  
268 *arXiv:1611.08024*, 2016.
- 269 [8] Mikhail A Lebedev and Miguel AL Nicolelis. Brain-machine interfaces: From basic science to neuroprostheses  
270 and neurorehabilitation. *Physiological reviews*, 97(2):767–837, 2017.
- 271 [9] Steven Lemm, Benjamin Blankertz, Thorsten Dickhaus, and Klaus-Robert Muller. Introduction to machine  
272 learning for brain imaging. *Neuroimage*, 56(2):387–399, 2011.
- 273 [10] Nanying Liang and Laurent Bougrain. Decoding finger flexion from band-specific ecog signals in humans.  
274 *Frontiers in neuroscience*, 6:91, 06 2012.
- 275 [11] Sergio Machado, Fernanda Arajo, Flavia Paes, Bruna Velasques, Mario Cunha, Henning Budde, Luis F Basile,  
276 Renato Anghinah, Oscar Arias-Carrin, Mauricio Cagy, et al. Eeg-based brain-computer interfaces: an overview  
277 of basic concepts and clinical applications in neurorehabilitation. *Reviews in the Neurosciences*, 21(6):451–468,  
278 2010.
- 279 [12] Joseph N Mak and Jonathan R Wolpaw. Clinical applications of brain-computer interfaces: current state and  
280 future prospects. *IEEE reviews in biomedical engineering*, 2:187–199, 2009.
- 281 [13] Yaron Meirovitch, Hila Harris, Eran Dayan, Amos Arieli, and Tamar Flash. Alpha and beta band event-related  
282 desynchronization reflects kinematic regularities. *Journal of Neuroscience*, 35(4):1627–1637, 2015.
- 283 [14] Luis Fernando Nicolas-Alonso and Jaime Gomez-Gil. Brain computer interfaces, a review. *Sensors*, 12(2):1211–  
284 1279, 2012.

- 285 [15] Miguel Pais-Vieira, Mikhail Lebedev, Carolina Kunicki, Jing Wang, and Miguel Nicolelis. A brain-to-brain  
286 interface for real-time sharing of sensorimotor information. *Scientific reports*, 3:1319, 02 2013.
- 287 [16] Johanna Reichert, Silvia Kober, Christa Neuper, and Guilherme Wood. Resting-state sensorimotor rhythm (smr)  
288 power predicts the ability to up-regulate smr in an eeg-instrumental conditioning paradigm. *Clinical Neurophys-*  
289 *iology*, 126, 02 2015.
- 290 [17] Yannick Roy, Hubert Banville, Isabela Maria Carneiro de Albuquerque, Alexandre Gramfort, and Jocelyn  
291 Faubert. Deep learning-based electroencephalography analysis: a systematic review, 01 2019.
- 292 [18] Gerwin Schalk and Eric C Leuthardt. Brain-computer interfaces using electrocorticographic signals. *IEEE*  
293 *reviews in biomedical engineering*, 4:140–154, 2011.
- 294 [19] Robin Tibor Schirrmeister, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katha-  
295 rina Eggenberger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with  
296 convolutional neural networks for brain mapping and decoding of movement-related information from the human  
297 eeg. *arXiv preprint arXiv:1703.05051*, 2017.
- 298 [20] Hyeyoung Shin, Robert Law, Shawn Tsutsui, Christopher Moore, and Stephanie Jones. The rate of transient beta  
299 frequency events predicts behavior across tasks and species. *eLife*, 6, 11 2017.
- 300 [21] Ksenia Volkova, Mikhail A Lebedev, Alexander Kaplan, and Alexei Ossadtchi. Decoding movement from elec-  
301 trocorticographic activity: A review. *Frontiers in neuroinformatics*, 13:74, 2019.
- 302 [22] Daniel Wolpert and Zoubin Ghahramani. Computational principles of movement neuroscience. *Nature neuro-*  
303 *science*, 3 Suppl:1212–7, 12 2000.
- 304 [23] Ivan Zubarev, Rasmus Zetter, Hanna-Leena Halme, and Lauri Parkkonen. Adaptive neural network classifier for  
305 decoding meg signals. *NeuroImage*, 197:425–434, 2019.