

---

# SPEECH DECODING FROM A SMALL SET OF SPATIALLY SEGREGATED MINIMALLY INVASIVE INTRACRANIAL EEG ELECTRODES WITH A COMPACT AND INTERPRETABLE NEURAL NETWORK

---

A PREPRINT

**Artur Petrosyan**  
Center for Bioelectric Interfaces  
Higher School of Economics  
Moscow, Russia

**Alexey Voskoboinikov**  
Center for Bioelectric Interfaces  
Higher School of Economics  
Moscow, Russia

**Dmitrii Sukhinin**  
Center for Bioelectric Interfaces  
Higher School of Economics  
Moscow, Russia

**Anna Makarova**  
Center for Bioelectric Interfaces  
Higher School of Economics  
Moscow, Russia

**Anastasia Skalnaya**  
Moscow State University of Medicine and Dentistry  
Moscow, Russia  
skalnaya95@gmail.com

**Nastasia Arkhipova**  
Moscow State University of Medicine and Dentistry  
Moscow, Russia  
exeast@gmail.com

**Mikhail Sinkin**  
Moscow State University of Medicine and Dentistry  
Scientific Research Institute of First Aid to them. N.V. Sklifosovsky  
Moscow, Russia mvsinkin@gmail.com

**Alexei Ossadtchi\***  
Center for Bioelectric Interfaces  
Higher School of Economics  
Artificial Intelligence Research Institute, AIRI  
Moscow, Russia ossadtchi@gmail.com

October 21, 2022

## ABSTRACT

1 **Background:** Speech decoding, one of the most intriguing BCI applications, opens up plentiful  
2 opportunities from rehabilitation of patients to direct and seamless communication between human  
3 species. Typical solutions rely on invasive recordings with a large number of distributed electrodes  
4 implanted through craniotomy. Here we explored the possibility of creating speech prosthesis in a  
5 minimally invasive setting with a small number of spatially segregated intracranial electrodes.

6 **Methods:** We collected one hour of data (from two sessions) in two patients implanted with invasive  
7 electrodes. We then used only the contacts that pertained to a single sEEG shaft or an ECoG stripe to  
8 decode neural activity into 26 words and one silence class. We employed a compact convolutional  
9 network-based architecture whose spatial and temporal filter weights allow for a physiologically  
10 plausible interpretation.

**Results:** We achieved on average 55% accuracy using only 6 channels of data recorded with a single minimally invasive sEEG electrode in the first patient and 70% accuracy using only 8 channels of data recorded for a single ECoG strip in the second patient in classifying 26+1 overtly pronounced words. Our compact architecture did not require the use of pre-engineered features, learned fast and resulted in a stable, interpretable and physiologically meaningful decision rule successfully operating over a contiguous dataset collected during a different time interval than that used for training. Spatial characteristics of the pivotal neuronal populations corroborate with active and passive speech mapping results and exhibit the inverse space-frequency relationship characteristic of neural activity. Compared to other architectures our compact solution performed on par or better than those recently featured in neural speech decoding literature.

**Conclusions:** We showcase the possibility of building a speech prosthesis with a small number of electrodes and based on a compact feature engineering free decoder derived from a small amount of training data.

## 1 Introduction

Brain-computer interfaces (BCIs) directly link the nervous system to external devices [20] or even other brains [44]. While there exist many applications of BCIs [1], clinically relevant solutions are of primary interest since they hold promise to rehabilitate patients with sensory, motor, and cognitive disabilities [35],[14].

BCIs can deal with a variety of neural signals [42, 33] such as, for example, electroencephalographic (EEG) potentials sampled with electrodes located on the scalp [34], or neural activity recorded invasively with intracortical electrodes penetrating cortex [26] or placed directly onto the cortical surface [49]. A promising and minimally invasive way to directly access cortical activity is to use stereotactic EEG (sEEG) electrodes inserted via a burr hole made in the skull. Recent advances in implantation techniques including the use of brain's 3D angiography, MRI and robot-assisted surgery help to further reduce the risks of such an implantation and make sEEG technology an ideal trade-off for BCI applications [23]. ECoG strips is another method to achieve direct electrical contact with cortical tissue with minimal discomfort to a patient [2].

The ability to communicate is vital to humans and speech is the most natural channel for it. Inability to speak dramatically affects the quality of life. A number of disorders can lead to a loss of this vital function, for example, cerebral palsy and stroke of the brain stem. Also, in some cases severe speech deficits may occur after a radical brain tissue removal surgery in oncology patients. While several technologies have been proposed to restore the communication function they primarily rely on brain controlled typing or imaginary handwriting [59] and appear to be practical only for severely affected patients. At the same time only in the United States 50 million people suffer from not being able to use their speech production machinery properly. A significant fraction of them have pathology not amenable by alaryngeal voice prosthesis [30] or "silent speech" devices [17] and require a neurally driven speech restoration solution.

Several successful attempts of BCI based speech restoration have already been made and a significant progress is achieved in decoding phonemes [60, 46, 40], individual words [36, 39, 55], continuous sentences [36, 39, 55] and even acoustic features [22, 55, 4] followed by the speech reconstruction algorithms using either Griffin-Lim or deep neural network algorithms inspired by WaveNet [4].

These solutions employ a broad variety of machine learning approaches for decoding speech from brain activity data. Starting from linear models [60], LDA [5], metric models [22] to deep neural networks (DNN) [36, 39, 55], that in general do not require manual feature engineering and can be applied directly to the data, however sometimes operating over a set of handcrafted features primarily derived from high-gamma activity. Several different neural network architectures have been tried for the speech decoding task: 1) relatively shallow ones consisting of a few convolutional or LSTM layers, 2) truly deep architectures with inception blocks [55] or with skip connections exploiting residual learning technique [4] as well as those borrowed from the computer vision applications [27, 56], 3) ensembles of DNN [39] making final solution more robust. Interestingly, that linear methods demonstrate comparable or, at least, close to DNNs decoding quality. Moreover, the latest studies obtained state of the art decoding accuracy using just a few layers over a set of handcrafted physiologically plausible features [36, 39]

The majority of the existing neural speech decoding studies rely on heavily multichannel brain activity measurements implemented with massive ECoG grids [39, 36, 4, 3] covering significant cortical area. These solutions for reading off brain activity are not intended for a long term use and are associated with significant risks to a patient [29] and suffer from a rapid loss of signal quality due to the leakage of the cerebrospinal fluid under the ECoG grid even if it is properly perforated. sEEG is a promising alternative whose implantation process is significantly less traumatic as compared to that of the large ECoG grids. The use of sEEG has already being explored for the speech decoding task

65 [5] but the reported decoder again relied on a high count of channels from multiple sEEG shafts distributed over a  
 66 large part of the left frontal and left superior temporal lobes which reduces the practicality of the proposed solution.  
 67 A solution capable of decoding speech from the locally sampled brain activity would be an important step towards  
 68 creating a speech prosthesis device.

69 Here we explore the possibility of decoding individual words from intracranially recorded brain activity sampled with  
 70 compact probes whose implantation did not require a full blown craniotomy. Our study comprises two subjects im-  
 71 planted either with sEEG shafts or ECoG stripes both via compact drill holes. We decode individual words using  
 72 either 6 channels of data recorded with a single sEEG shaft or the 8 channels sampled using a single ECoG strip. For  
 73 decoding we employed our interpretable CNN architecture [45] augmented with the bidirectional LSTM layer [25] to  
 74 compactly model local temporal dependencies in the internal speech representation that we used as the intermediate  
 75 decoding target. We also compared the ultimate word decoding accuracy achieved with different internal representa-  
 76 tions. Our decoder operated causally using only the data from time intervals preceding the decoded time moment and  
 77 therefore is fully applicable in a real-time decoding setting. Overall our study is the first attempt to achieve accept-  
 78 able individual words decoding accuracy from cortical activity sampled with compact non-intracortical probes whose  
 79 implantation is not likely to cause significant discomfort to a patient and can be done even with local anesthesia.

## 80 2 Data

81 In this study we used two datasets collected from two epilepsy patients undergoing planned sEEG and ECoG implan-  
 82 tation for the needs of presurgical mapping. The first patient was implanted bilaterally with a total of 5 sEEG shafts  
 83 with 6 contacts in each with the goal to localize seizure onset zone. The implantation was performed under general  
 84 anesthesia via five 3-mm drill holes. The second patient was implanted with 9 ECoG stripes of 8 contacts each cover-  
 85 ing frontal and inferior temporal lobes. The implantation was performed via several 12 mm drill holes. Figure 1  
 86 demonstrates post-surgical CT scans of the two patients. On the second day past the implantation both patients went  
 87 through the active and passive [51] speech mapping procedures that yielded concordant results. In Patient 1 electrical  
 88 stimulation of the 10-11 pair (300  $\mu$ s, 2.5 mA, 50 Hz) resulted in pronounced speech arrest. The passive speech  
 89 mapping procedure based on computing the mutual information (MI) between the speech envelope and the envelope  
 90 of the gamma-band (60 Hz -100 Hz) filtered sEEG activity resulted into a sharp peak of the MI values for electrodes  
 91 9-12, see Figure 1.a.c. No speech related artifacts were observed when stimulating contacts 11-12 which could be due  
 92 to the very sparing stimulation settings used in this patient - our stimulation current in this patient never exceeded 3  
 93 mA which is below the traditionally average current magnitude typically used for speech mapping [15]. Stimulation  
 94 based speech mapping in Patient 2 caused involuntary tongue retraction when applied between electrodes 15-16 and  
 95 the MI profile highlighted contacts 13-15, see Figure 1.b,d. Note that the exact shape of the MI profiles depends on  
 96 the filtering parameters and therefore these plots need to be interpreted carefully. The MI profiles may also confuse  
 97 speech production and one’s own speech perception processes, especially given the observation demonstrated in [32]  
 98 that gamma activity in the auditory cortex accurately tracks the perceived speech envelope and may contribute to the  
 99 observed MI.

100 The study was conducted according to the ethical standards of the 1964 Declaration of Helsinki. The participants  
 101 provided written informed consent prior to the experiments. The ethics research committee of the National Research  
 102 University, The Higher School of Economics approved the experimental protocol of this study. After the patients  
 103 signed the appropriate informed consent we asked them to read off the succession of the 6 sentences presented at a  
 104 comfortable pace in randomized order on the computer screen. Each sentence was repeated 30 and 65 times by the  
 105 first and the second patient respectively. The sentences contained on average 4.3 words. Half of the sentences had  
 106 direct and the other half indirect order of words and the majority of words within a single sentence started from the  
 107 same letter. This was done to enable subsequent neurolinguistic analysis of the collected datasets. The sEEG in Patient  
 108 1 was recorded with an 80 channel g.HIamp amplifier. Patient’s 2 ECoG was registered with a 64 channel EBNeuro  
 109 BE Plus LTM device. The sampling rate was set to  $F_s = 19200$  Hz (Patient 1) and  $F_s = 4096$  Hz (Patient 2). In both  
 110 cases synchronously with neural activity we recorded speech signal measured with Behringer XM8500 microphone.

## 111 3 Methods

### 112 3.1 Data preprocessing

113 We first parsed audio data into separate words. To this end, we manually processed several example word alignments  
 114 and then used them to find similar ones by means of the dynamic time warping (DTW) algorithm [8]. Manual check  
 115 shows that the absolute majority of the word alignments were detected correctly. For each word this procedure resulted  
 116 in a list of index pairs corresponding to the start and the end of the word’s utterance. Audio data were processed using

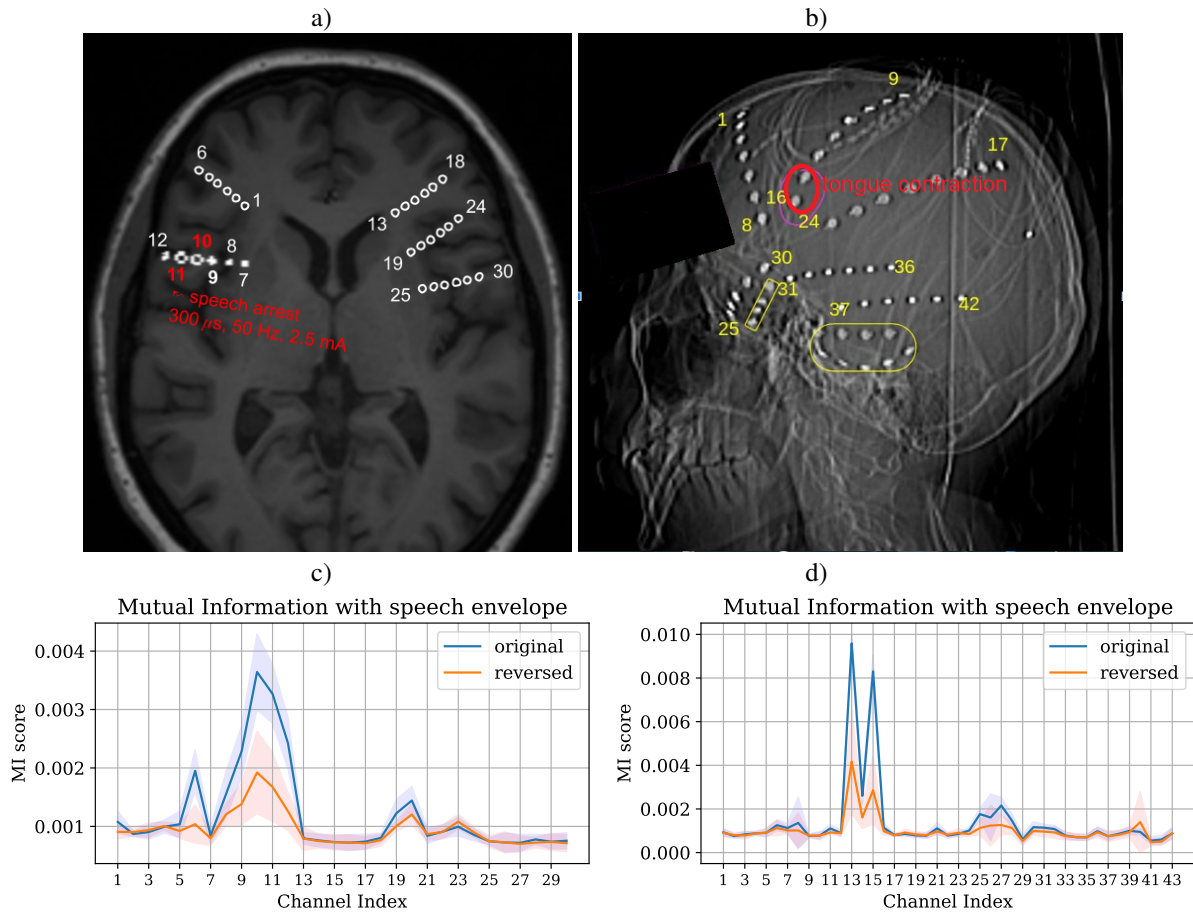


Figure 1: Figure 1: a) sEEG contacts extracted from the post-implantation CT scan of the first patient superimposed over her MRI. Bipolar electrical stimulation of the 10-11 pair (300  $\mu$ s, 2.5 mA, 50 Hz) resulted in reproducible speech arrest. b) CT of the second patient who was implanted with nine 8-contact ECoG stripes covering bilateral frontal and inferior temporal lobes. Bipolar electrical stimulation applied to electrodes 15-16 caused involuntary tongue retraction. c) Patient 1, mutual information profile between the speech envelope and gamma-band (60 Hz -100 Hz) filtered sEEG activity shows a sharp peak of the MI values for electrodes 10 and 11. d) Patient 2, mutual information profile between the speech envelope and gamma-band (70 Hz -100 Hz) shows peak over electrodes 13-15 and in several other locations of this ECoG strip. The MI profiles between the time-reversed audio stream and original ECoG data is shown in red. The shadow corresponds to the standard deviation of the MI values estimated using the collection of different 3 minute long segments. The remaining bumps in the time-reversed MI profile may be due to the inherently rhythmic nature of the audio stream produced by the patient in response to the sequence of computer instructions.

117 Librosa Python software package [38] in order to extract log-mel spectral coefficients (LMSC) [54], mel-frequency  
 118 cepstral coefficients (MFCC) [61] and several derivatives of the linear predictive coding (LPC) coefficients as de-  
 119 scribed below. The sequences of these internal speech representations (ISRs) were downsampled to 1 KHz. We have  
 120 experimented with all listed ISRs, see section 3.3 and Figures 7, 11.b in section 4. In the majority of the reported  
 121 results we used LMSCs as the internal speech representation.

122 sEEG and ECoG data went through minimal preprocessing such as causal band-pass FIR filter in the 5-150 Hz fre-  
 123 quency range. Then the data were resampled to 1 kHz sampling rate and the amplitude in each channel was standard-  
 124 ized by subtracting the mean and dividing by the standard deviation. This multichannel data was used as an input of  
 125 our decoding algorithm.

## 126 3.2 Decoding

127 In accord with the view expressed in [39] we decided to explore decoding accuracy on the level of individual words.  
 128 On the one hand, words represent a sufficiently low level of detail which permits extension of the obtained solution  
 129 into a broader range of application scenarios. On the other hand, words are less volatile as compared to phonemes  
 130 as the articulation of the latter greatly varies depending on the flanker sounds neighboring the phoneme. This may  
 131 mean that the neural encoding governing the transition between the different states of the articulatory tract may vary  
 132 significantly from case to case depending on the phonetic context a phoneme is encountered in.

## 133 3.3 Internal speech representations

134 Most of the ISRs are based on modeling speech signal as produced by an excitation sequence passing through a linear  
 135 time-varying filter [28]. The excitation sequence is the air flow in the larynx and the filter is formed by the articulatory  
 136 tract elements (pharynx, vocal folds, tongue, lips, teeth) whose mutual geometry changes over time.

137 Linear predictive coding (LPC) and cepstral analysis are the two principal ways to estimate parameters of such a filter.  
 138 LPC analysis is based on a direct estimate of the auto-regressive model coefficients  $a_i$  through Burgs method [37].  
 139 However, these prediction coefficients themselves are unstable, as their small changes may lead to large variations in  
 140 the spectrum and possibly unstable filters. In order to decrease such an instability the following several equivalent  
 141 representations are commonly used.

142 Reflection coefficients (RC)  $k_i$  can be computed alongside with prediction coefficients through Burgs method and  
 143 represent the ratio of the amplitudes of the acoustic wave reflected by and the wave passed through a discontinuity.

144 Another descriptor, log-area ratio (LAR) coefficients,  $g_i$ , are equal to the natural logarithm of the ratio of the areas of  
 145 adjacent sections in a lossless tube equivalent of the vocal tract having the same transfer function and can be computed  
 146 from the reflection coefficients as  $g_i = \ln \left( \frac{1-k_i}{1+k_i} \right)$ .

147 Line spectral frequencies (LSF) is another highly efficient speech data compression technique [52] as errors in repre-  
 148 senting one coefficient generally result in a spectral change only around that frequency.

149 In what follows we will present our experiments with several ISRs but our final decoding accuracy results are based  
 150 on the use of log-mel spectral coefficients (LMSC).

### 151 3.3.1 Synchronous decoding

152 Our goal is to decode specific words from the immediately preceding chunks of neural activity data. The direct  
 153 approach would require gathering a large amount of training data. Instead we developed our decoding solution based  
 154 on the idea described in [36] where the vocoder-like compact ISR is used for regularization purposes during the  
 155 training. However, here instead of using the ISR as a regularizer we employ it as the intermediate target. In other  
 156 words, we first use our compact and interpretable architecture [45] to decode the ISR vector (e.g.  $M = 40$  LMSCs)  
 157 from either sEEG or ECoG based measurements of brain activity. After having trained this ISR decoder optimizing the  
 158 average correlation coefficient between the actual and the decoded ISRs we fix its weights and train a convolutional  
 159 neural network to decode discrete words based on the representations that emerged in one before the last layer of ISR  
 160 decoder network. After training, our two-stage architecture operates as a single network on the minimally preprocessed  
 161 neural activity data and yields discrete classification of individual words at its output.

162 For the training word classification task we semi-automatically, see section 3.1, extracted alignment of each word.  
 163 We used only a chunk of neural data that corresponds to the particular word's alignment (we do not use information  
 164 of neighborhood words). We also added a "silent" class that corresponds to the intervals of silence between word



189 extraction parameters and let our architecture learn them during the training process guided by the optimization of the  
 190 mean (across all ISR elements) Pearson’s correlation coefficient between the original and the decoded ISR timeseries.  
 191 This feature extraction is performed by the adaptive envelope detector (ED) block that comprises a succession of the  
 192 factorized spatial and temporal convolution operations followed by the rectification and smoothing blocks. The ED  
 193 during training can potentially adapt to extracting instantaneous power of specific neuronal populations activity pivotal  
 194 for the downstream task of predicting the ISRs. In the search for the optimum, the ED weights are not only tuned to  
 195 such a target source but also tuned away from the interfering sources [21, 45]. The proper interpretation of the learnt  
 196 ED’s weights allows for subsequent discovery of the target source geometric and dynamical properties.

197 In order to measure the quality, we used 6 fold cross validation. In each of the folds we took 5/6 of independent  
 198 sessions and 1/6 for validation, corresponding to 50 minutes and 10 minutes respectively. We used Adam optimiser  
 199 for training with  $\alpha = 0.0003$  learning rate parameter. For training we used the entire train portion of our data, not just  
 200 speech segments. Using only speech intervals worsened the quality of decoding by 15-20%. Our intuition here is that  
 201 non-speech segments are also useful to the training and may serve regularization purposes. Also, hypothetically, the  
 202 brain activity that determines the upcoming utterance happens during the silence interval and therefore not including  
 203 this segment into the training could have detrimental effect on the final classification accuracy.

204 In the majority of our experiments we used LMSCs as the ISR, but as described in section 4 we have also experimented  
 205 with the other ISRs outlined in section 3.3, as a target for our first network. After having trained our compact architec-  
 206 ture to decode the ISRs as our intermediate target we used a 2D-convolution ResNet to perform discrete classification  
 207 of 26 words and the silent class using the representations developed in the one before the last layer of the compact  
 208 architecture, see Figure 3.

209 Importantly, our experiments show that the use of the internal representations that emerged in the LSTM layer instead  
 210 of the actual decoded ISRs noticeably improves the final word classification accuracy. This observation is inline with a  
 211 similar finding in a completely different domain [19] where the authors advocated the use of multiple separate "views"  
 212 generated by different networks as the input to the final classifier in the image classification task.

### 213 3.5 Performance metrics

214 We use correlation coefficient to measure ISR–from–neural activity reconstruction quality. To assess the words decod-  
 215 ing accuracy when operating in the synchronous mode we downsample "silence" intervals to avoid the positive bias  
 216 in the reported numbers and then measure accuracy as the fraction of correctly classified utterances. We report our  
 217 results in the form of  $27 \times 27$  confusion matrices illustrating the proportion of correct and erroneous decoding of 26  
 218 words and the silence class.

219 To assess the accuracy when operating in the asynchronous mode we use precision-recall characteristics. As described  
 220 earlier, see Figure 2, for each  $i$ -th word we compute smoothed probability profiles  $\tilde{p}_i(t)$  for each time instance  $t$ . We  
 221 make a decision about a word being pronounced only at time points corresponding to the local maximums of  $\tilde{p}_i(t)$  that  
 222 cross the threshold  $\theta$ . The  $i$ -th word is decoded if the local maximum of  $\tilde{p}_i(t)$  located above  $\theta$  also appears to be the  
 223 largest among all other profiles, i.e.  $\tilde{p}_k(t)$ ,  $k \neq i$ .

224 In case the chosen  $i$ -th word (or the silence) corresponds to the one that is currently being uttered we mark this event  
 225 as true positive (TP). If after such a detection  $\tilde{p}_i(t)$  remains above the threshold and exhibits another local maximum  
 226 which exceeds the values of all other smoothed probability profiles we will also make a decision to "utter" the  $i$ -th  
 227 word. However, in this case this decision will be marked as false positive (FP) even if  $t$  belongs to the time range  
 228 corresponding to the actual  $i$ -th word, because this results in the duplicated uttering and adds errors to the decoded  
 229 words sequence. We also mark as FP the events when the index of the detected word does not match that of the actual  
 230 pronounced word, see Figure 13.a for the graphical representation of the above description. To compute these PR  
 231 curves we first smooth the probability profiles delivered by the neural network with a simple box-car averaging over  
 232 the 0.2 sec segment. Then, we vary the detection threshold (single value for the entire test data segment) and compute  
 233 the corresponding precision-recall pair. Doing so for a dense grid of thresholds we obtain a threshold independent  
 234 metrics of algorithms performance.

235 We represent our asynchronous decoding results in the form of precision-recall curves parameterised by the threshold  
 236  $\theta$  applied to probability profiles. Since our decoder uses softmax at its output we smoothly varied threshold  $\theta$  in  $(0, 1)$   
 237 range to calculate *precision* and *recall* indicators for each value of the threshold according to the expressions:

$$precision = \frac{TP}{TP + FP}, \quad recall = \frac{TP}{N}, \quad (1)$$

238 where  $N$  is the total number of actual utterances performed by the patient.

239 The obtained curves characterize the amount of information present at the decoder output and facilitate comparison of  
 240 various solutions. In practice, however, when such an asynchronous BCI is used by a patient the specific value of the  
 241 threshold is to be set based on the user’s preferences.

### 242 3.6 Weights interpretation

243 When dealing with overt speech decoding from neural activity data one needs to make sure that the obtained decision  
 244 rule indeed uses neural activity data and does not exploit for decoding the possible artifacts such as electrical currents  
 245 accompanying muscular activity or the acoustic signal leaked into neural data channels via, for example, microphone  
 246 effect [47]. Thankfully, the widely spread over cortex spatial patterns of muscular activity occupying high frequency  
 247 range [16] can be delineated from neural signals whose high frequency components, on the contrary, tend to be  
 248 restricted to spatially compact cortical regions [57, 41]. To do so one needs access to both spatial and frequency  
 249 domain patterns of the activity that appears pivotal to the decoder. Interpretable decision rules facilitate such tests  
 250 for physiological plausibility of the obtained solutions. By extracting spatial and frequency domain patterns from the  
 251 weights of the corresponding layers [21, 45] we can check for the physiological plausibility using domain specific  
 252 knowledge as described above.

253 In this work we use our compact convolutional network as the front-end which allows for the theoretically justified  
 254 interpretation of its spatial and temporal convolution weights by extracting spatial and frequency domain patterns  
 255 corresponding to the neuronal populations whose activity is pivotal to the specific downstream task. The details of our  
 256 approach are outlined in [45], next we briefly review the basic ideas behind it.

257 The front-end of our network comprises factorized spatial and temporal convolution layers, see Figure 3. During  
 258 training, the spatial and temporal filter weights of each branch not only get tuned to the pivotal neuronal sources but  
 259 also tune away from the interfering signals.

260 In terms of spatial processing, that is combining the data from different sensors with specific weights, each branch of  
 261 our adaptive envelope detector (ED), see Figure 3, corresponds to the model studied in [21]. However, each branch  
 262 of the ED contains both spatial and temporal filters. Therefore, as we show in [45], the interpretation of branch’s  
 263 *spatial* weights needs to be conducted within the context set by the corresponding *temporal* filter. Since both spatial  
 264 and temporal filtering are linear, interchanging them in the above statement is also valid and thus branch’s temporal  
 265 filter weights interpretation needs to be done taking into account the spatial filter of this branch. More formally our  
 266 approach is summarized below and in Figure 3.

267 Our ED processes data in chunks of a prespecified length of  $N$  samples. First, assume that the input segment length  
 268 is equal to the filter length in the 1-D temporal convolution layer. Consider a chunk of input data from  $L$  channels  
 269 observed over the interval of  $N$  time moments that can be represented by matrix  $\mathbf{X}[n] = [\mathbf{x}[n], \mathbf{x}[n-1], \dots, \mathbf{x}[n-N+1]] \in \mathbb{R}^{L \times N}$ . Processing of  $\mathbf{X}[n]$  by the first two layers performing spatial and temporal filtering can be described for  
 270 the  $m$ -th branch by a bi-linear product as  
 271

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

272 where  $\mathbf{w}_m \in \mathbb{R}^L$  is a vector of spatial weights and  $\mathbf{h}_m \in \mathbb{R}^N$  is a vector temporal weights for branch  $m$ . The non-  
 273 linearity,  $\text{ReLU}(-1)$ , in combination with the low-pass filtering performed by the second convolutional layer (that  
 274 smooths the rectifier output  $r_m[n]$ ) and extracts the envelopes  $e_m[n]$  of the rhythmic signals.

275 We assume that upon training the spatial unmixing coefficients and temporal filter impulse responses implement op-  
 276 timal processing and tune each branch of our architecture to a specific neuronal population with its characteristic  
 277 geometric and dynamical properties. But it is crucial to realize that under Wiener optimal condition each branch not  
 278 only gets tuned to a specific population but also tunes itself away from the interfering activity. As detailed in [45],  
 279 assuming that channel timeseries are zero-mean random processes the underlying neuronal population topographies  
 280 can be found as

$$\mathbf{g}_m = \mathbb{E}\{\mathbf{y}_m[n] \mathbf{y}_m^T[n]\} \mathbf{w}_m^* = \mathbf{R}_m^y \mathbf{w}_m^* \quad (3)$$

281 where  $\mathbf{R}_m^y = \mathbb{E}\{\mathbf{y}_m[n] \mathbf{y}_m^T[n]\}$  is a  $L \times L$  spatial covariance matrix of the temporally filtered data  $\mathbf{y}_m[n] = \mathbf{X}[n] \mathbf{h}_m$ ,  
 282  $L$  is the number of input channels. Thus, when interpreting individual spatial weights corresponding to each of the  
 283  $M$  branches of the architecture shown in Figure 3 one has to take into account the temporal filter weights  $\mathbf{h}_m^*$  of this  
 284  $m$ -th branch.

285 The temporal weights should be interpreted in a similar way, i.e. taking into account the corresponding spatial fil-  
 286 ter. Assuming that channel timeseries are zero-mean random processes,  $N$  is the number of taps in the temporal



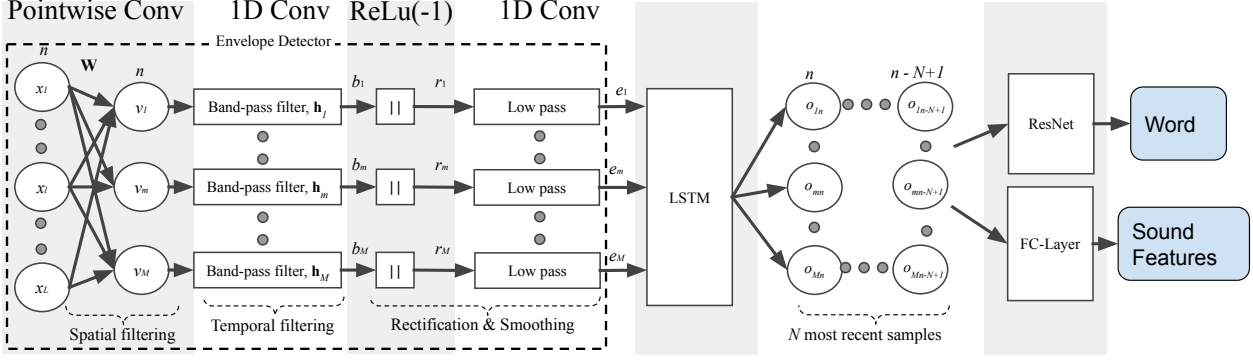


Figure 3: The architecture based on [45] and adapted for speech classification task. We used the same envelope detector technique to extract robust and meaningful features from the neuronal data. We then used the LSTM layer to account for the sequential structure of the speech ISR (e.g. LMSC) and finally decoded it with a fully connected layer over the LSTM hidden state ( $o_{ij}$  on the figure). A separate 2D convolutional network was trained and used to classify separate words from the activity of this pretrained LSTM.

287 convolution filter  $\mathbf{h}_m^*$ , the temporal pattern is given by

$$\mathbf{q}_m = \mathbb{E}\{\mathbf{v}_m[n]\mathbf{v}_m^T[n]\}\mathbf{h}_m^* = \mathbf{R}_m^v \mathbf{h}_m^* \quad (4)$$

288 where  $\mathbf{R}_m^v = \mathbb{E}\{\mathbf{v}_m[n]\mathbf{v}_m^T[n]\}$  is an  $N \times N$  tap covariance matrix of an  $N$ -samples long chunk of spatially filtered  
 289 data  $\mathbf{v}_m[n] = \mathbf{X}[n]^T \mathbf{w}_m = [v_m[n], v_m[n-1], \dots, v_m[n-N+1]]^T$ .

290 As shown in [45] if we relax the assumption about the length of the data chunk being equal to the length of the temporal  
 291 convolution filter we can arrive at Fourier domain representation of the second-order dynamics of the neuronal popula-  
 292 tion the  $m$ -th branch is tuned to. The power spectral density  $Q_m(f)$  of this population's activity can be derived from  
 293 the power spectral density (PSD)  $P_{v_m}(f)$  of the spatially filtered input data  $v_m[n]$  and the Fourier transform  $H_m(f)$   
 294 of the temporal weights vector  $\mathbf{h}_m(f)$  as in (5):

$$Q_m(f) = P_{v_m}(f)H_m(f) \quad (5)$$

295 The important distinction that contrasts our weights interpretation approach from the methodology used in the majority  
 296 of reports utilizing neural networks with separable spatial and temporal filtering operations is that our procedure  
 297 accounts for the fact that during training the spatial filter formation is taking place within the context set by the  
 298 corresponding temporal filter, and vice versa. Also, in [45] the authors for the first time introduced the notion of  
 299 the frequency domain pattern  $Q_m(f)$  of neuronal population's activity. Note that  $Q_m(f)$  vs.  $H_m(f)$  has the same  
 300 difference as the spatial pattern vs. spatial filter weights which was brilliantly illustrated earlier in [21].

301 Using the expressions 3 and 5 we can explore the corresponding spatial and frequency domain patterns of each trained  
 302 branch (head) of our decoding architecture. If our architecture latched to the data of neuronal origin then the spatial  
 303 patterns of larger extent should correspond to sources with frequency domain patterns localized to lower frequency  
 304 ranges and vice versa. Such mutual relationship if observed may reassure that our decoder relies on genuinely neuronal  
 305 information.

## 306 4 Results

### 307 4.1 Microphone effect

308 To exclude the possibility of data leak associated with electric contacts capacitance change driven by the acoustic  
 309 speech signal vibes, also known as microphone effect [48], we compared spectral content of the recorded neural data  
 310 and that of the speech signal in 0-2000 Hz frequency range. Time-frequency diagrams corresponding to a typical  
 311 20 seconds long segment of a representative channel of neuronal data and the acoustic signal are shown in Figure 6  
 312 for two patients. Visual analysis does not reveal the characteristic banded structure of speech signal (lower row) in  
 313 the time-frequency profiles of the neuronal data (top row). To perform an objective assessment for all channels of  
 314 neural data we calculated the correlation coefficient between the temporal profiles of the instantaneous power in each  
 315 frequency band of neural and acoustic data. We have then used a permutation test to assess statistical significance of the  
 316 observed correlation coefficients to be non-zero. To this end we split the acoustic data into segments corresponding to

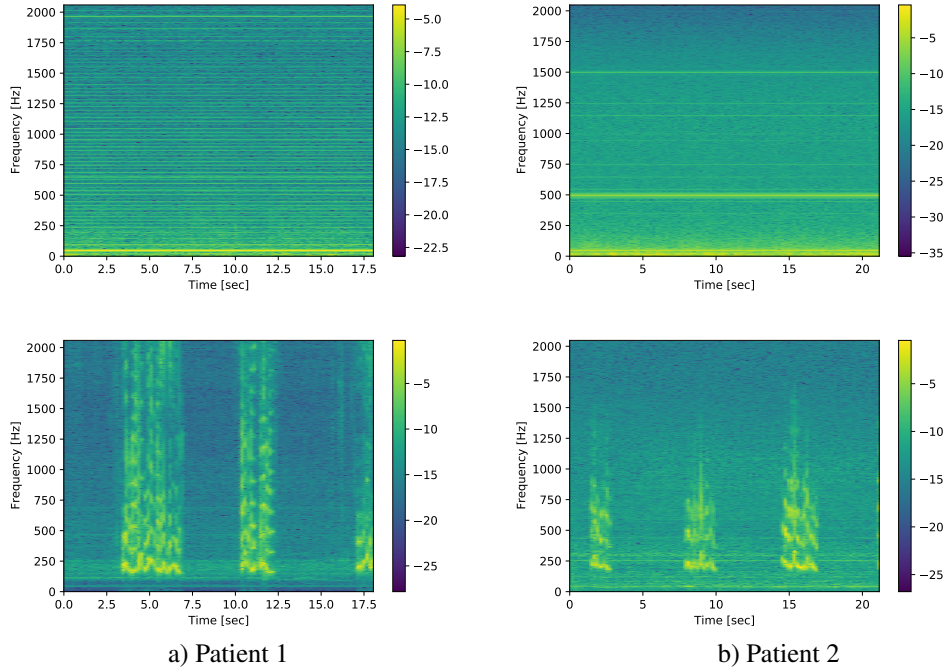


Figure 4: Log-spectrograms for audio and neural data in the electrode with the strongest audio-neural correlations

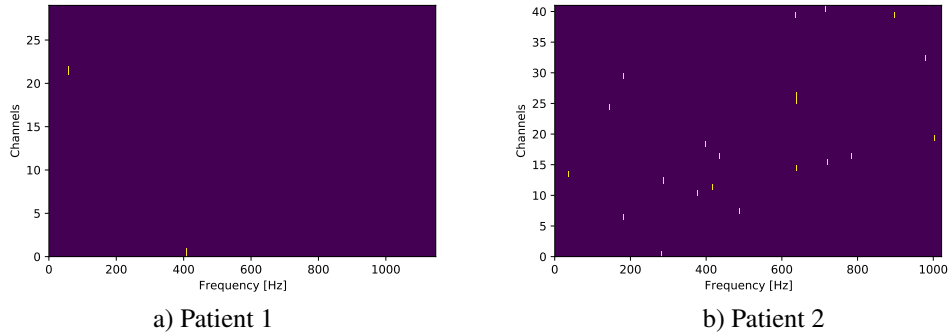


Figure 5: Pairs (channel, frequency) with rejected  $H_0$  hypotheses that correlation is zero at  $\alpha = 0.05$ . See more description in the corresponding section.

317 word utterances and randomly shuffled 10000 times the order of such segments to destroy the original correspondence  
 318 between the neuronal and acoustic data in order to compute surrogate correlation coefficient distribution for each  
 319 (channel, frequency) pair. Then we have computed the asymptotic  $p$ -values as the fraction of times when the surrogate  
 320 correlation coefficients appeared to be greater than the correlation coefficients observed in the original non-shuffled  
 321 data. To correct for multiple comparisons due to running a massive set of tests for all (channel, frequency) pairs we  
 322 used the BH FDR correction procedure [9] and obtained a set of adjusted  $p$ -values. Those (channel, frequency) pairs  
 323 whose corresponding adjusted  $p$ -values fall below 0.05 are highlighted in Figure 5 and do not show any systematic  
 324 segregation in neither of the two patients.

325 The above analysis assures that there were no identifiable effects of acoustic information leakage into the data channels  
 326 carrying neural activity signals.

## 327 4.2 Decoding internal speech representation

328 In this study we mainly focused on the contacts confined to a single stereo-EEG shaft in Patient 1 or a single stripe  
 329 in Patient 2. To select the specific contiguous block of contacts we have computed mutual information between the

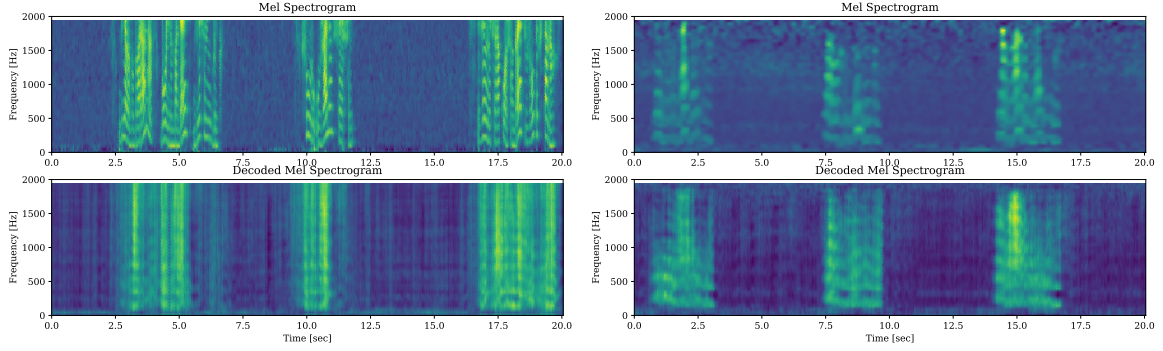


Figure 6: Example of a true and decoded from neural activity log mel-spectrograms.

330 speech envelope and the envelope of the high-gamma band cortical activity signal for each channel, see Figure 1 c)  
 331 and d) panels. We have observed a clear delineation in the amount of mutual information between different electrodes.  
 332 Reassuringly, high MI values closely matched electrodes whose stimulation led to speech arrest in Patient 1 and tongue  
 333 contraction in Patient 2, see Figure 1. Red curves represent the amount of MI computed using the mutually reversed  
 334 neuronal and audio sequences.

335 Some remaining values of the MI in the reversed sequence can be explained by the rhythmic structure of the computer  
 336 instructions to utter that the patients followed. Although the MI profiles are sensitive to the filtering option, see section  
 337 2, we still consider it a useful tool for delineation between task related and unrelated channels. As we show in Figure  
 338 11 exploiting the MI informed selection of channel groups yields the best decoding accuracy that matches the value  
 339 achieved with the entire set of channels.

340 As evident from Figure 7 our compact architecture using only 6 sEEG channels from a single sEEG shaft achieved  
 341 about 65% mean correlation over  $M = 40$  LMSCs in Patient 1 and almost 60% for Patient 2 with 8 channels from a  
 342 single ECoG stripe. These accuracy values in decoding internal speech representation are comparable to those reported  
 343 in [4] where significantly greater count of data channels collected by multiple sEEG shafts was used. An example of  
 344 the original and decoded 40 LMSCs is shown in Figure 6 for two patients.

345 We have also experimented with decoding several other internal speech representations (ISRs) as shown in the left  
 346 panel of Figure 7. Each color corresponds to a specific ISR method. For both patients we display the ISRs using the  
 347 same order. Interestingly, in both patients LMSCs appeared to be decoded best, PCs followed and got closely matched  
 348 by the MFCCs. The reflection coefficients had the worst decoding accuracy. As we will show next, however, this order  
 349 is not retained when we use words classification accuracy as a criterion. Most likely the mean correlation coefficient  
 350 between the true and decoded ISRs is determined by their specific statistical properties and the extent to which the  
 351 fluctuations in their coefficients reflect changes between the silence and speech intervals. To explore this we have  
 352 computed masked correlation coefficients using only the intervals when the actual speech was present. As expected  
 353 the mean correlation coefficient dropped significantly and the order in which the different ISRs lined up changed as  
 354 well, see the right panel of Figure 7. LMSCs on average still remained among the ISRs with top decoding accuracy  
 355 followed by the LPC coefficients and MFCC.

356 Each ISR is a vector and instead of the average values shown in Figure 7 in Figure 8 we present the decoding accuracy  
 357 values achieved for each of the elements in the three ISRs with the best average decoding accuracy: LMSCs, MFCCs  
 358 and LPC coefficients. Here we also observe similar tendencies for both patients. For each we show the histograms of  
 359 correlation coefficients computed over the entire data range (blue) and only over speech intervals (orange).

360 The achieved so far ISR decoding accuracy does not yield intelligible speech when, for example, the recovered LMSC  
 361 sequence is converted back into the sound. Nevertheless, as we will show next the decoded LMSC profiles and other  
 362 ISRs support the classification of discrete words sufficiently well.

### 363 4.3 Words decoding in synchronous mode

364 We achieved 55% accuracy using only 6 channels of data recorded with a single minimally invasive sEEG electrode in  
 365 the first patient in classifying 26+1 overtly pronounced words (3.7% chance level). The left panel of Figure 9 shows  
 366 the corresponding confusion matrix and the individual decoding accuracy values for each word in Patient 1.

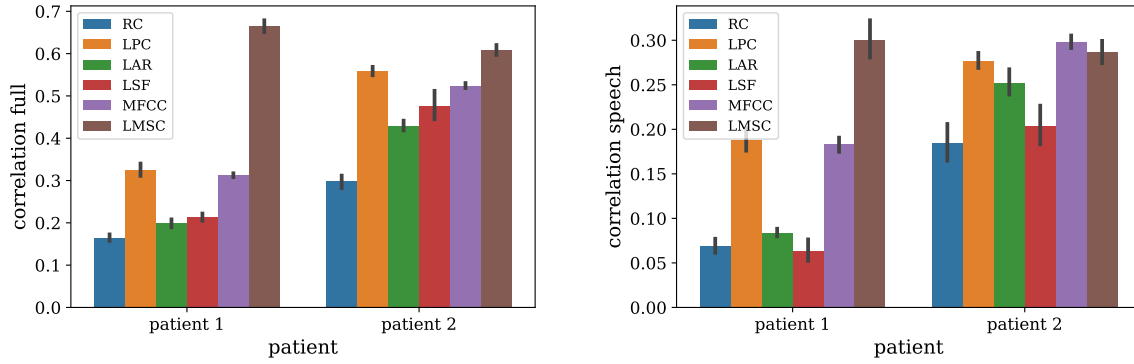


Figure 7: Comparison of the decoding accuracy achieved for different ISRs: LPC - linear predictive coding coefficients, LSF - line spectral frequencies, RC - reflection coefficients, LAR - log-area ratios, LMSCs - log-mel spectrograms, MFCC - mel-frequency cepstral coefficients. To test for statistical significance of the observed differences in decoding quality, we performed Wilcoxon signed-rank tests with Bonferroni correction: (\*) - p-value is less than 0.05, (\*\*) - 0.01, (\*\*\*) - 0.001. We added this information to the caption. The left panel corresponds to the correlation coefficients between the actual and decoded temporal profiles computed over the entire time range of the test data segment. Statistically significant differences for: Patient 1 - LMSC with RC/LPC/LAR/LSF/MFCC (\*\*\*), RC with LPC/LSF/MFCC (\*\*\*), RC with LAR (\*), LPC with LAR/LSF (\*\*\*), MFCC with LAR/LSF (\*\*\*). Patient 2 - RC with LPC/LAR/LSF/MFCC/LMSC (\*\*\*), LAR with LPC/MFCC/LMSC (\*\*\*), MFCC with LMSC (\*\*\*), LMSC with LPC/LSF (\*\*), LPC with MFCC (\*). In the right panel the correlation coefficient is computed only over the time intervals where the actual speech was present. Statistically significant differences for: Patient 1 - LMSC with RC/LPC/LAR/LSF/MFCC (\*\*\*), RC with LPC/MFCC (\*\*\*), LPC with LAR/LSF (\*\*\*), MFCC with LAR/LSF (\*\*\*). Patient 2 - RC with LPC/MFCC/LMSC (\*\*), LSF with MFCC/LMSC (\*\*), LAR with RC/MFCC (\*), LPC with LSF (\*).

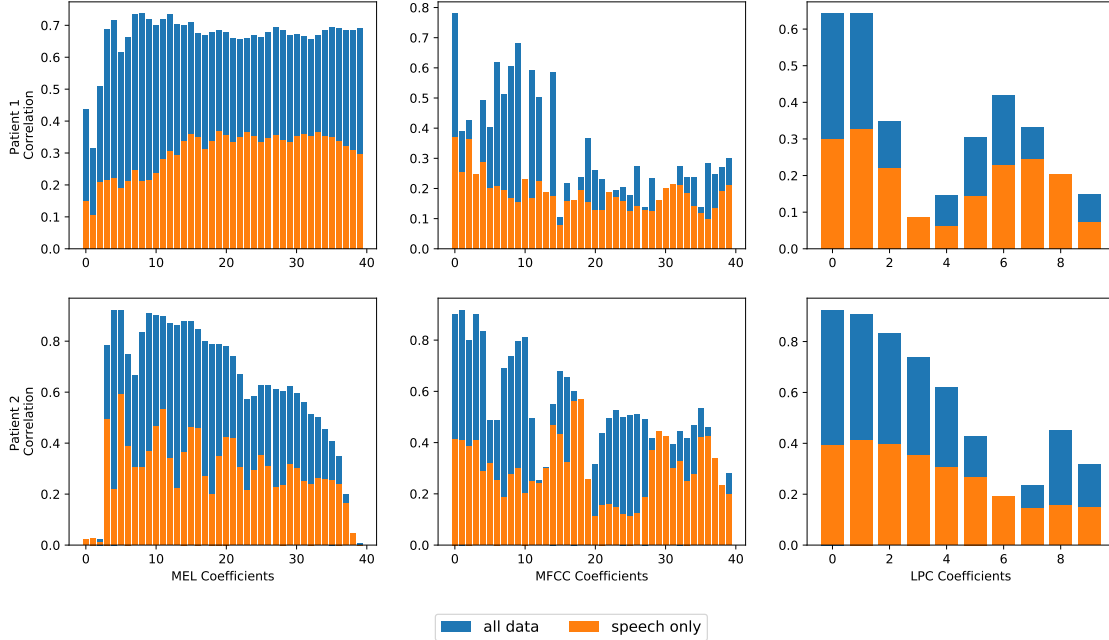


Figure 8: Correlation coefficient of predicted and actual ISR elements for two patients (rows). The two overlaid histograms correspond to the correlations computed over the entire time range (blue) and only over the speech intervals (orange).

367 Spatial characteristics of the first three branches corresponding to the most pivotal neuronal population are shown in  
 368 the left column of Figure 10.a. We can see that dominantly the activity of these pivotal populations is mapped onto

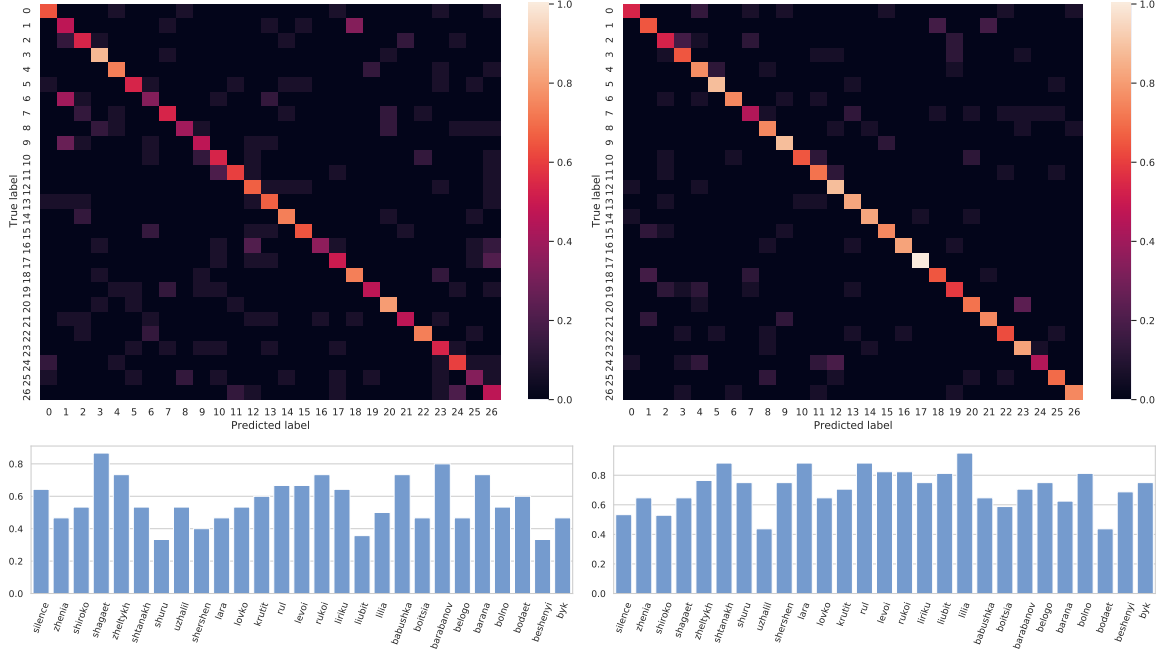


Figure 9: Confusion matrix of classified words for patient 1 and patient 2. Words list: 0. silence, 1. zhenia, 2. shiroko, 3. shagaet, 4. zhelytkh, 5. shtanakh, 6. shuru, 7. uzhalil, 8. shershen, 9. lara, 10. lovko, 11. krutit, 12. rul, 13. levoi, 14. rukoi, 15. liriku, 16. liubit, 17. lilia, 18. babushka, 19. boitsia, 20. barabanov, 21. belogo, 22. barana, 23. bolno, 24. bodaet, 25. beshenyi, 26. byk. In the bottom we show the individual word decoding accuracy values, corresponding to the diagonal of the confusion matrix

369 electrodes with indices 9 to 12. This corroborates with the results of an active speech mapping procedure where we  
 370 found that bipolar electrical stimulation of electrodes indexed 10 and 11 resulted in transient speech arrest as shown  
 371 in Figure 1 a). Frequency domain patterns presented next to the corresponding spatial patterns illustrate physiological  
 372 plausibility. First of all the top branch has activity not only in the lower frequency range but also in the traditional  
 373 gamma band and this branch corresponds to the spatially compact pattern highlighting a single channel with index 12.  
 374 At the same time, the two branches are characterized by frequency domain patterns concentrated over relatively lower  
 375 frequency range. Interestingly, and in agreement with [57] these branches have relatively more spread out spatial  
 376 patterns as compared to that of the first branch.

377 Similar analysis is shown for Patient 2 in the right panel of Figure 9 and Figure 10.b. In this patient implanted with  
 378 ECoG grids we have achieved on average 70% of words decoding accuracy. We can also observe a striking trend  
 379 where spatially more compact populations are characterized by the activity in the higher frequency bands.

### 380 4.3.1 Weights interpretation

381 The advantage of the DNN based approach is that it does not require manual feature engineering, however, these  
 382 methods are typically over-parameterized and exhibit greedy behaviour. Such a greediness in the neurophysiological  
 383 context may result in the network latching on signals of non-neuronal origin. This problem can be monitored in  
 384 compact, domain knowledge driven architectures equipped with a proper weights interpretation approach. To this end  
 385 we have applied the recently developed approach detailed in [45]. Our goal here is to explore the mutual relation  
 386 between the spatial and frequency domain patterns each branch of our compact DNN architecture got tuned to during  
 387 the training process. This analysis will also help us to exclude the fact that our network exploits muscular activity  
 388 associated with the speech production process. The principles behind this analysis have been briefly outlined in  
 389 section 3.6.

390 The result of applying our weights interpretation procedure to each of the three branches of our compact DNN is shown  
 391 in Figure 10. We illustrate both spatial (left column) and frequency domain (right column) patterns of the neuronal  
 392 populations for each of the three most significant branches of our network. The frequency domain plot also contains



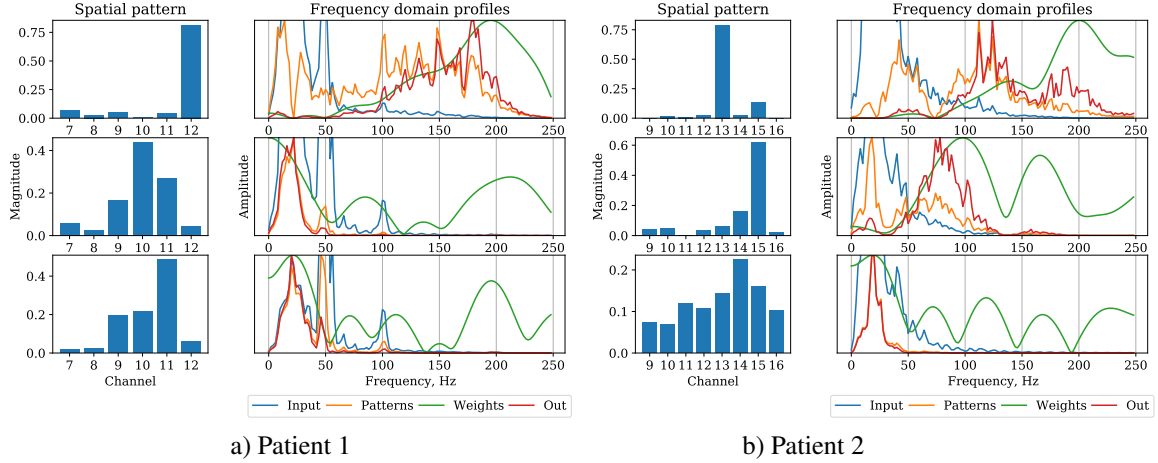


Figure 10: Theoretically justified weights interpretation applied to the most significant branches of the architecture in Figure 3. Orange trace in the top panel shows the power spectral density pattern of the activity of the neuronal population this branch is tuned to. The left panel shows the spatial pattern of this population. We can conclude that this source dominantly projects onto the 12th contact located at the lateral part of the sEEG electrode (shaft), see Figure 1 a). Similar picture is observed for the second patient present in the two right columns. Here we also observe a striking trend of spatially more compact populations being characterized by the activity in the higher frequency bands.

393 the curves corresponding to the power spectral density (PSD) of the input timeseries obtained by the spatial filtering  
 394 of the multichannel data at the input of the network and the PSD of the branch’s output timeseries.

395 From the top row of patterns corresponding to the first branch of the decision rule for Patient 1 we can see that the  
 396 PSD occupies a high 100-200 Hz frequency range and the corresponding spatial pattern is confined to only a single  
 397 channel with index 12. At the same time the second branch with a much more spread out spatial pattern occupying  
 398 channels 9-11 is characterized by the PSD confined to the lower 10-40 Hz frequency range. The reciprocal space-  
 399 frequency relation that hallmarks neuronal activity and distinguishes it from the electro-muscular artifacts is also very  
 400 well pronounced in the second patient. Moving downwards we observe the gradual growth of the spatial spread with  
 401 the PSD frequency range migrating from the higher to lower frequency range.

402 Combined together with domain knowledge [12, 13, 11, 57] highlighting reciprocal space-time relationship in the  
 403 observed cortical activity patterns and phenomenological observations [16] on the properties of the electromuscular  
 404 activity and its representation on the cortex the observed combinations of the spatial and PSD patterns allow us to  
 405 make a conclusion regarding the neuronal origin of the data our decoder latched on during the training process. The  
 406 analysis for microphone effect reported in section 4.1 also excludes the possibility that the decoding is done based on  
 407 the acoustic signal leaking into neuronal data channels.

408 In this patient we have witnessed certain discrepancy between the stimulation based mapping and the electrode indexes  
 409 that resulted from our weights interpretation procedure where electrode 12 was highlighted, yet speech production  
 410 problems were registered when stimulating contacts 10-11, but see Figure 1 a,c. This could have resulted from the  
 411 very sparing stimulation settings used in this patient - our stimulation current never exceeded 3 mA which is below the  
 412 traditional average current magnitude typically used for speech mapping [15].

413 For Patient 2, weights interpretation of the three most important branches of our network show the primary involvement  
 414 of electrodes 13, 14 and 15 into the decoding process which is partially congruent with stimulation based speech  
 415 mapping, see Figure 1.b,d, where we found that the stimulation applied between 15-16th electrodes yielded reliable  
 416 tongue retraction (back from the requested tongue protrusion state). In this patient we also observe a very pronounced  
 417 reciprocity in the space-frequency patterns. As shown in the right panel of Figure 10 moving from the top to the  
 418 bottom we observe how a very compact spatial pattern transitions into a more spatially spread out one. At the same  
 419 time, the corresponding frequency domain patterns tend to move leftwards so that the most compact spatial pattern  
 420 corresponds to the activity with the highest central frequency. This is the expected property of neural activity that has  
 421 been highlighted earlier in several studies [41, 57].

422 Approximate MNI coordinates of electrodes found to be pivotal for decoding in both patients are given in Table 1.  
 423 In Patient 1 stereo-EEG electrodes with the majority of contacts located deep in the sulci of the left operculum. The  
 424 location corresponds to Brocas region whose activity is traditionally registered in a broad range of language related

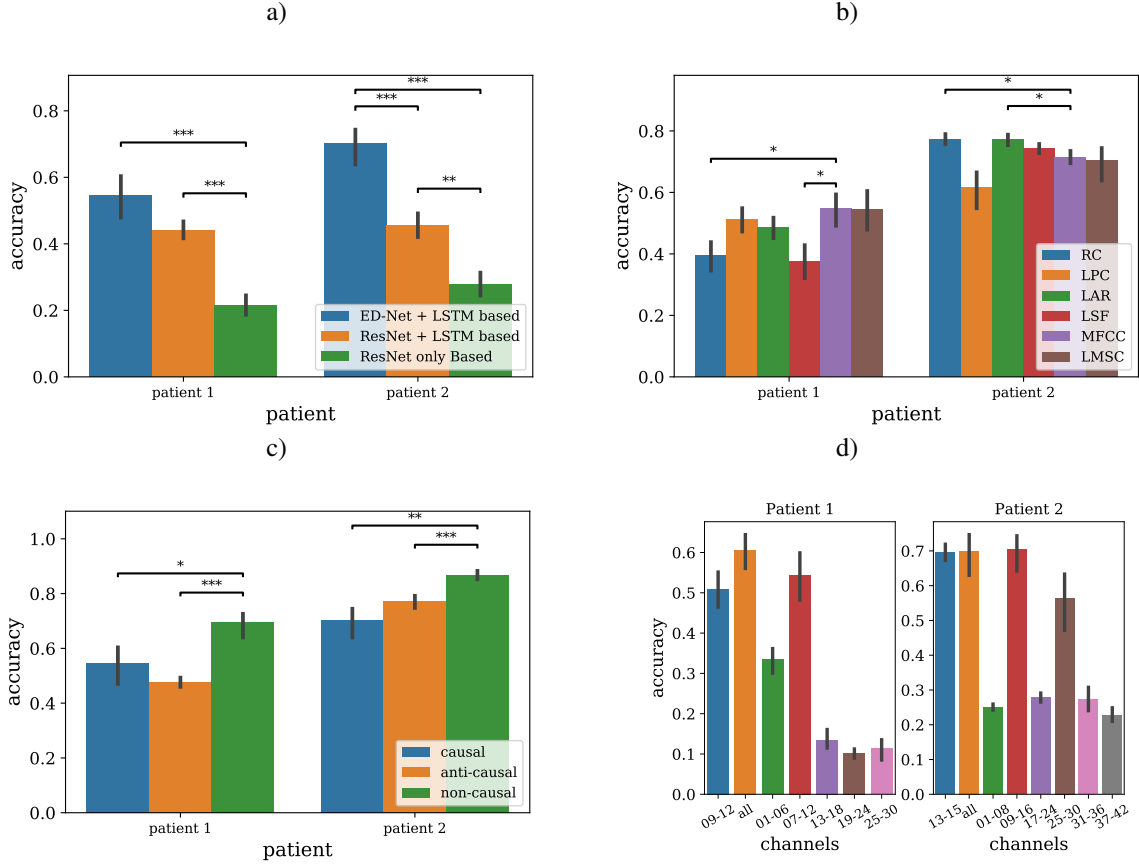


Figure 11: Comparative analysis. To test for statistical significance of the observed differences in ISR reconstruction fidelity, we performed Wilcoxon signed-rank test with Bonferroni correction: (\*) - p-value is less than 0.05, (\*\*) - 0.01, (\*\*\*) - 0.001. We added this information to the graph. In cases the brackets appeared to overload the plot we placed the statistical testing results in this caption. a) Comparison of different neural network models. b) Comparison of different possible intermediate sound representation, LPC - linear predictive coding coefficients, LSF - Line Spectral Frequencies, RC - reflection coefficients, LAR - log-area ratios, LMSC - log-mel spectrogram coefficients, MFCC - mel-frequency cepstral coefficients. c) Comparison of different possible lag. d) Comparison of decoding quality for different subset of channels. Statistically significant differences for: Patient 1: 09-12 with 01-06/13-18/19-24/25-30 (\*\*\*) , 09-12 with all (\*) , all with 01-06/13-18/19-24/25-30 (\*\*\*) , 01-06 with 13-18/19-24/25-30 (\*\*\*) , 07-12 with 13-18/19-24/25-30 , 07-12 with 01-06 (\*\*) . Patient 2: 13-15 with 01-08/17-24/31-36/37-42 (\*\*\*) , all with 01-08/17-24/31-36/37-42 (\*\*\*) , 09-16 with 01-08/17-24/31-36/37-42 (\*\*\*) , 17-24 with 37-42 (\*\*\*) , 17-24 with 01-08 (\*) , 25-30 with 01-08/31-36/37-42 (\*\*\*) , 25-30 with 17-24 (\*\*\*) , 25-30 with 13-15/all/09-16 (\*)

425 tasks. For patient 2 the pivotal electrodes cover the inferior portion of the precentral gyrus. Our stimulation results in  
 426 Patient 2 do not quite match the anatomical location and functionally better correspond to the ventral precentral gyrus,  
 427 the structure located inferior to the precentral gyrus and known to house the tongue motor area. This could be due  
 428 to atypical organization of the cortex in this patient. Locations of electrodes in both patients are remote with respect  
 429 to the belt area (MNI: -58, -28, 13) whose gamma-band activity was shown to reliably track the perceived speech  
 430 envelope (Kubanek et al., 2013). These functional and anatomical arguments together with the causal approach to the  
 431 ISR decoding reduce the chance that our decoder operation is based on the subjects own speech perception.

432 In the above we have analyzed spatial and frequency domain patterns of the neuronal populations that were found to  
 433 be pivotal to the ISR decoding task and forming the internal representations to be subsequently used as an input to our  
 434 words classification network.

435 For the front-end network weights interpretation to make sense in the context of the word classification task we also  
 436 need to demonstrate the dependence of the final word classification accuracy on the fidelity of the individual ISR  
 437 decoding achieved by the front-end network. To this end we have performed additional experiments. We rerun the  
 438 training and terminated it at different points to yield various ISR decoding accuracy and then subsequently trained

Patient 1, stereo-EEG electrode		Patient 2, ECoG strip	
Electrode index	MNI Coordinates	Electrode index	MNI Coordinates
9	-40,19,4	13	-65,-24,43
10	-45,20,4	14	-66,-18,39
11	-50,20,4	15	-67,-13,34
12	-53,21,4	25	-35,17,-46
		30	-56, 8,-21

Table 1: MNI coordinates of pivotal electrodes

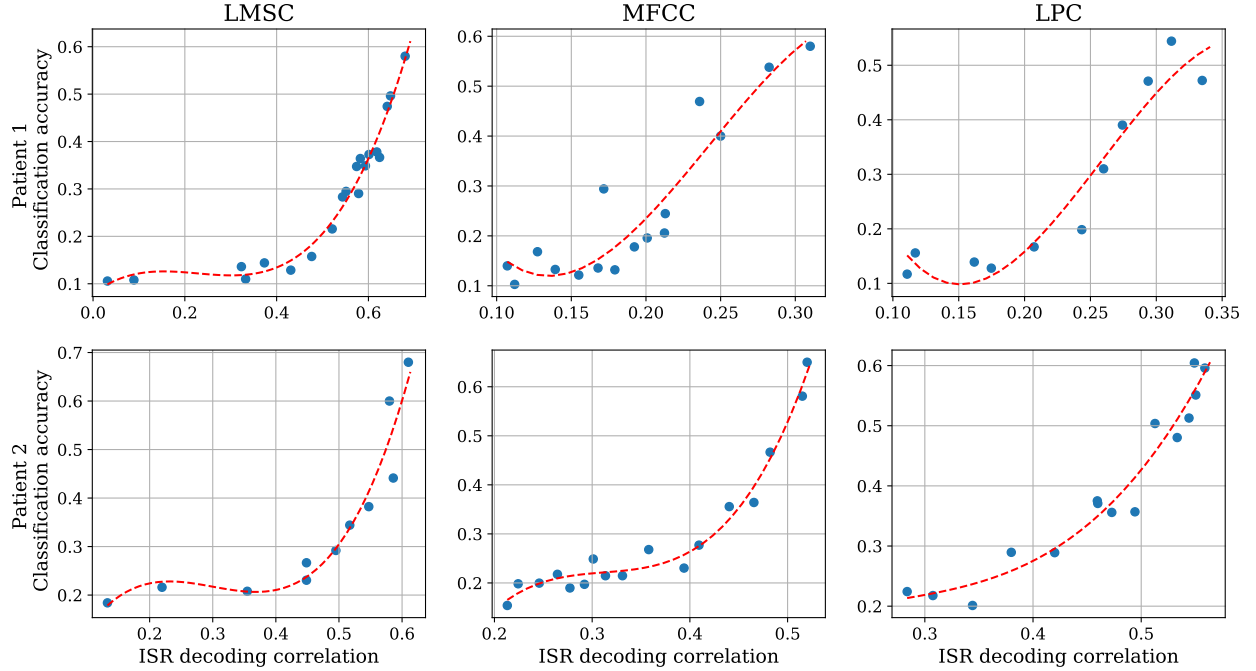


Figure 12: Dependence of the final word classification accuracy on the decoded vs. true ISR correlation. Red line is the third order trend fitted to the data to facilitate visual perception.

439 our word classification network. In Figure 12 we show the observed dependence of the words classification accuracy  
 440 on the correlation coefficient between the actual and the decoded ISRs (LMSC, MFCC and LPC). Indeed for each  
 441 ISR we witness the direct relation between its decoding fidelity as measured by the correlation coefficient and the  
 442 corresponding discrete words classification accuracy.

#### 443 4.4 Comparative analysis

444 In this work we employed the compact architecture, see Figure 3, that comprises multiple branches of envelope detec-  
 445 tors (ED) of spatially filtered data whose output is fed into the LSTM layer followed by a fully connected network.  
 446 This architecture uses factorized spatial and temporal filters that get adapted during training and allows for interpreta-  
 447 tion of the filter weights into the spatial and spectral patterns as demonstrated in Figure 10. These patterns can then be  
 448 used to infer location and dynamical properties of the underlying neuronal populations.

449 Here we compared this network to several other architectures. We found that out of several neural networks only  
 450 Resnet-18 offers a comparable, although significantly worse, performance when used instead of the ED block in our  
 451 architecture, see Figure 3. The LSTM layer also appears to be very useful in capturing the dynamics of features  
 452 extracted either with ED or ResNet blocks, see Figure 11.a. We hypothesize that this situation may be caused by the  
 453 adequate balance in the number of parameters to be tuned for the ED-based network and the amount of data available  
 454 for training as compared to several other more sophisticated architectures.



455 Words decoding accuracy results reported in Figure 9 correspond to the case when 40 LMSCs were used to train the  
 456 front-end ISR decoder network, see Figure 3. We have also experimented with several other ISRs as described in  
 457 section 3.3 and presented the results in Figure 11.b.

458 Interestingly, the differences in the individual ISR decoding fidelity, see Figure 7, does not transfer into the corre-  
 459 sponding words classification accuracy where all of the ISRs yield more or less comparable performance. A possible  
 460 explanation here could be that some ISRs in addition to the information regarding the sequence of the articulatory tract  
 461 configurations (corresponding to a specific sequence of phonemes and invariant to the pitch, timbre, loudness, etc.)  
 462 contain the information about purely acoustic features of the utterance such as fundamental frequency, voice timbre,  
 463 local volume, etc., which could be easier to decode than the articulatory tract parameters critical for the words classifi-  
 464 cation task. The subsequent words classification largely requires only the first type of information and therefore may  
 465 yield comparable words classification performance for the different ISRs as long as all of them contain this essential  
 466 information.

467 The reported ISR and words decoding accuracy results are presented for the causal processing mode, i.e. when the  
 468 data window strictly precedes the time-point the prediction is made for. We have also experimented with anti-causal  
 469 (the window is strictly in the future w.r.t. to the predicted time-point) and non-causal (when the data window covers  
 470 pre- and post- intervals around the point in question). These results are plotted in Figure 11.c. In both patients we see  
 471 the best performance when the data-window is allowed to be located both in the future and in the past w.r.t. the point  
 472 to be predicted. This result is expected since in the non-causal setting the algorithm can use information about the  
 473 cortical activity that occurs in response to the uttered word.

474 In this work we mainly focused on decoding from a small number of contacts confined either to a single stereo-EEG  
 475 shaft or an ECoG stripe. In both cases the electrodes can be implanted without a full-blown craniotomy via a drill hole  
 476 in the skull. We have chosen the particular subset of contacts using the mutual information (MI) metric, see Figure  
 477 1 which closely matched stimulation-based mapping results. Both of our patients were implanted with several sEEG  
 478 shafts or ECoG stripes, see Figure 1. In Figure 11.d we show the results of a similar analysis but using other subsets  
 479 of electrodes located on the other shafts or stripes. Noteworthy is that MI based selection yielded significantly better  
 480 performance as compared to the other spatially segregated electrode groups.

481 According to Figure 1.d electrodes 25-27 also show the increased MI values between the ECoG and acoustic envelope.  
 482 This corroborates with the results in Figure 11.d where the use of the stripe with these electrodes yielded the second  
 483 best decoding accuracy. The stripe is placed in the inferior region of the left anterior temporal cortex and the MNI  
 484 coordinates of the first (25) and the last (30) electrodes from this stripe are given in Table 1. According to [53] these  
 485 areas appear to be active during the implicit comprehension of spoken and written language. Given that the sentences  
 486 we used slightly deviate from the standard sentences used in daily life and are likely to require some additional effort  
 487 and very mild emotional response beyond just mechanical reading. According to Figure 1 of [24], our electrodes  
 488 25-30 fall in the area 6e that appears to host representations of emotional words, see their Table 2. Finally, based on  
 489 [10], the temporal pole region where electrodes 25-30 are placed could be a part of the network that links temporal  
 490 pole with posterior structures to support thematic semantic processing during language production. When interpreting  
 491 these results we can not discount the mounting evidence that speech production and comprehension share neural  
 492 representation and speech production processes are not only localized to the left hemisphere but also involve bilaterally  
 493 distributed linguistic network [50] which explains advanced decoding accuracy in the speech decoding setting reliant  
 494 on bilaterally distributed electrodes [23].

#### 495 4.5 Asynchronous decoding of words

496 Traditionally, BCI can be used in two different settings: synchronous and asynchronous. In the synchronous setting  
 497 a command is to be issued within a specific time window. Usually, a synchronous BCI user is prompted at the  
 498 start of such a time window and has to produce a command (alter his or her brain state) within a specified time frame.  
 499 Therefore, the decoding algorithm is aware of the specific segment of data to process in order to extract the information  
 500 about the command. In the asynchronous mode the BCI needs to not only decipher the command but also determine  
 501 the fact that the command is actually being issued. The delineation between synchronous and asynchronous modes is  
 502 most clearly pronounced in BCIs with discrete commands implying the use of a categorical decoder.

503 In BCIs that decode a continuous variable, e.g. hand kinematics, such delineation between synchronous and asyn-  
 504 chronous modes is less clear. The first part of our BCI implements a continuous decoder of the internal speech  
 505 representation (ISR) features. Should this decoding appear of sufficient accuracy it could have been simply used as  
 506 an input to a voice synthesis engine. Such a scenario has already been implemented in several reports [4, 3] but these  
 507 solutions use a large number of electrodes which may explain better quality of ISR decoding. In our setting we aimed  
 508 at building a decoder operating with a small number of ecologically implanted electrodes and decided to focus on

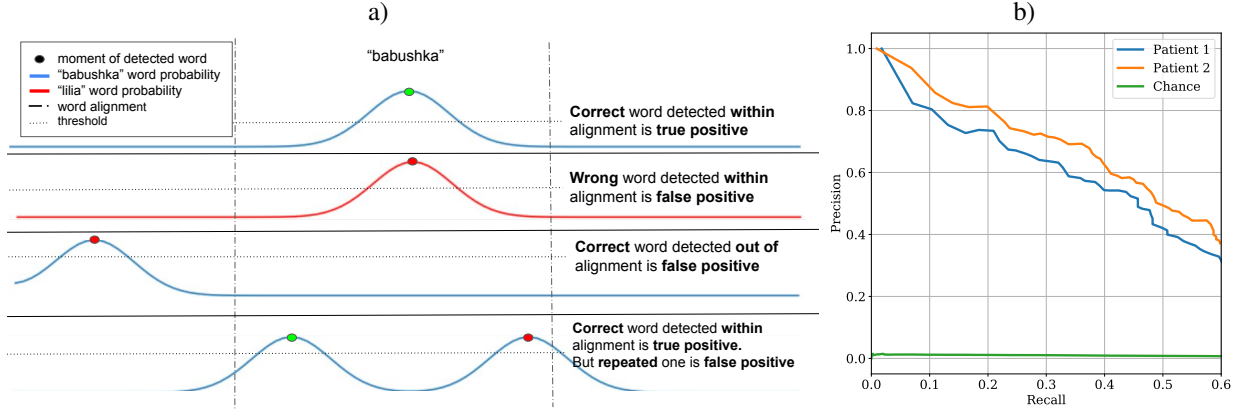


Figure 13: a) For each  $i$ -th word we compute smoothed probability profiles  $\tilde{p}_i(t)$  for each time instance  $t$ . The decision is then made about a word being pronounced only at time points corresponding to the local maximums of  $\tilde{p}_i(t)$  that cross the threshold  $\theta$ . In case the chosen  $i$ -th word matches the one that is currently being uttered we mark this event as true positive (TP). If after such a detection  $\tilde{p}_i(t)$  remains above the threshold and exhibits another local maximum which exceeds the values of all other smoothed probability profiles the  $i$ -th word is "uttered" again, but this event is marked as false positive (FP) even if  $t$  belongs to the time range corresponding to the actual  $i$ -th word. b) PR curves for asynchronous words decoding task. As in regular binary classification problem, in order to get PR curves we vary the detection threshold from 0 to 1 and for every fixed threshold value we compute the corresponding precision-recall pair. The detection threshold affects how many words will be «uttered» by our algorithm. Low detection threshold "utter" a lot of words and lead to high recall and low precision. And vice versa, high detection threshold "utter" only high confident words and lead to low recall and high precision. Note that definition of precision and recall is slightly different from conventional binary classification PR curves (see equation 1, figure 13.a and section 3.5 for details). We also show a chance level PR curve.

509 decoding individual words. We first used the continuously decoded ISRs to classify 26 discrete words and one silence  
 510 state in the synchronous manner. To implement this we cut the decoded ISR timeseries around each word's utterance  
 511 and use them as data samples for our classification engine.

512 To gain insight into the ability of our BCI to operate in a fully asynchronous mode we performed the additional analysis  
 513 as described in section 3.3.2. Figure 13 .b illustrates the performance of our BCI operating in a fully asynchronous  
 514 mode when the decoder is running over the succession of overlapping time windows of continuously decoded ISRs  
 515 and the decision about the specific word being uttered is made for each of such windows, see Figure 2. To quantify the  
 516 performance of our asynchronous speech decoder we used precision-recall curves as detailed in section 3.5 and Figure  
 517 13.a.

518 Although the observed performance significantly exceeds the chance level, it is not yet sufficient for building a full  
 519 blown asynchronous speech interface operating using a small number of minimally invasive electrodes. In our view and  
 520 based on the experience with motor interfaces, specific protocols to train the patient including those with immediate  
 521 feedback to the user [6] are likely to significantly improve the decoding accuracy in such systems which will boost the  
 522 overall feasibility of minimally invasive speech prosthetic solutions.

## 523 5 Conclusion

524 We have explored the possibility of building a practically feasible speech prosthesis solution operating on the basis of  
 525 neural activity recorded with a small set of minimally invasive electrodes. Implantation of such electrode systems does  
 526 not require a full craniotomy and combined with algorithmic solutions equipped with a joint human-machine training  
 527 protocol may form a basis for the future minimally invasive speech prosthesis.

528 There exist several reports exploiting intracortical activity recorded with Utah array like systems for speech prosthesis  
 529 purposes [60, 58, 18]. These recordings give access to the activity of individual neurons but remain potentially harmful  
 530 to the cortical tissue. In contrast, stentodes [43], electrodes located inside blood vessels and implanted using stent  
 531 technology, offer a potentially plausible solution for obtaining high quality brain activity signals without any kind of  
 532 craniotomy. These electrodes, however, unlike the intracortical arrays, register the superposition of neuronal activity  
 533 stemming from a large number of neuronal populations. Also, unlike the ECoG grids used in the majority of speech

534 prosthesis research these stent electrodes are confined to a relatively small volume. The signals measured in our setting  
535 with a small number of spatially confined sEEG and ECoG contacts can be considered as a proxy of the data collected  
536 by the stent electrodes and the signal processing approaches developed here could be potentially applied to stent electrode data  
537 in order to pave the road towards craniotomy-free speech BCI solutions.

538 We build our decoder using a two-step procedure. First we construct an interpretable architecture to decode the  
539 continuous internal speech representation (ISR) profiles from the neural activity and fix the weights of this compact  
540 neural network. In this case the particular ISR (LMSC, MFCC, LPC coefficients) is merely a target to train this  
541 front-end network. Then, when applying this network to neural activity data we take its hidden state before the last  
542 fully connected layer and use its activation as an input to the discrete classifier to distinguish between neural activity  
543 patterns corresponding to 26 words and one silence state. This approach resembles [36]. However, based on our  
544 experiments we found that replacing concurrent training of two classifiers with such a two step process improved the  
545 achieved decoding accuracy in our setting.

546 We have also paid particular attention to interpreting the obtained decision rule. Our main concern here was to exclude  
547 the possibility of using non-neural activity patterns in the overt speech decoding setting. To do so we exploited the  
548 concept of spatial and frequency domain patterns that pertain to the neuronal populations that each of the branches  
549 of our front-end network got tuned to. Several reports [16, 31, 7] explored the spatial and frequency domain patterns  
550 that manifest muscular activity in the subdural space. These are typically hallmarked with high-frequency spectra  
551 and large spatial extent which is the opposite to neural activity where we expect higher frequency activity to be more  
552 spatially confined as compared to the signals in the lower frequency bands. We applied the methodology described in  
553 [45] to recover spatial and frequency patterns of the underlying pivotal activity and found that they well adhered to the  
554 described properties of neural activity. We also did not find any evidence of microphone effect [47] in our data.

555 The accuracy we obtained in the synchronous mode appears sufficient to make a system usable in a real-life scenario  
556 when each word is "uttered" within a specific time slot, starting, for example, with a beep prompt. The extent to which  
557 the observed accuracy is transferred to a patient who lacks the ability to speak greatly depends on the specific medical  
558 case. Although we explored various arrangements of the data time window around the decision point our main results  
559 correspond to the decoder operating causally, i.e. utilizing neural activity strictly from the past which is expected not  
560 to depend on the perceived speech, see also [32]. This ensures that the observed accuracy can potentially transfer to  
561 real patients with speech function deficits given the appropriate patient training tools are developed.

562 Asynchronous BCI setting is clearly a more natural one for speech prosthesis operation. We experimented with our  
563 decoder in this scenario and observed a reasonable performance which however, needs to be improved before it can  
564 be used in practice. We recall 40% moments when one of the 26 words is uttered and in 60% of cases we correctly  
565 guessed this word out of 26 possible alternatives.

566 The use of a language model is known to improve speech decoding accuracy [55] and can also be added to improve the  
567 performance of the final consumer solution. However, our goal here was to assess to which extent the neural activity  
568 alone can be informative with regard to individual words classification and therefore we have deliberately refrained  
569 from using any language model in this study.

570 Overall our study showcases the possibility of building speech prosthesis with a small number of electrodes and based  
571 on a compact feature engineering free decoder derived from several tens of minutes worth of training data. To be  
572 translated into clinical practice this solution needs to be augmented with patient training procedures and a methodology  
573 to non-invasively determine implantation sites that would yield the best speech decoding accuracy.

## 574 **Acknowledgment**

575 This work is supported by the Center for Bioelectric Interfaces NRU HSE, RF Government grant, AG. No. 075-15-  
576 2021-624

577 **References**

- 578 [1] Sarah N Abdulkader, Ayman Atia, and Mostafa-Sami M Mostafa. Brain computer interfacing: Applications and  
579 challenges. *Egyptian Informatics Journal*, 16(2):213–230, 2015.
- 580 [2] Abidemi B Ajiboye and Robert F Kirsch. Invasive brain–computer interfaces for functional restoration. In  
581 *Neuromodulation*, pages 379–391. Elsevier, 2018.
- 582 [3] Hassan Akbari, Bahar Khalighinejad, Jose L Herrero, Ashesh D Mehta, and Nima Mesgarani. Towards recon-  
583 structing intelligible speech from the human auditory cortex. *Scientific reports*, 9(1):1–12, 2019.
- 584 [4] Miguel Angrick, Christian Herff, Emily Mugler, Matthew C Tate, Marc W Slutzky, Dean J Krusienski, and Tanja  
585 Schultz. Speech synthesis from ecog using densely connected 3d convolutional neural networks. *Journal of*  
586 *neural engineering*, 16(3):036019, 2019.
- 587 [5] Miguel Angrick, Maarten Ottenhoff, Lorenz Diener, Darius Ivucic, Gabriel Ivucic, Sofoklis Goulis, Jeremy Saal,  
588 Albert J Colon, Louis Wagner, Dean J Krusienski, et al. Real-time synthesis of imagined speech processes from  
589 minimally invasive recordings of neural activity. *bioRxiv*, 2020.
- 590 [6] Miguel Angrick, Maarten Ottenhoff, Lorenz Diener, Darius Ivucic, Gabriel Ivucic, Sophocles Goulis, Albert J  
591 Colon, Louis Wagner, Dean J Krusienski, Pieter L Kubben, et al. Towards closed-loop speech synthesis from  
592 stereotactic ecog: A unit selection approach. In *ICASSP 2022-2022 IEEE International Conference on Acoustics,*  
593 *Speech and Signal Processing (ICASSP)*, pages 1296–1300. IEEE, 2022.
- 594 [7] Tonio Ball, Markus Kern, Isabella Mutschler, Ad Aertsen, and Andreas Schulze-Bonhage. Signal quality of  
595 simultaneously recorded invasive and non-invasive eeg. *Neuroimage*, 46(3):708–716, 2009.
- 596 [8] Kalaba Bellman. Bellman r., kalaba r. *On adaptive control processes*, *IRE Trans. Autom. Control*, 4(2):1–9,  
597 1959.
- 598 [9] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to  
599 multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- 600 [10] Deena Schwen Blackett, Jesse Varkey, Janina Wilmskoetter, Rebecca Roth, Keeghan Andrews, Natalie Busby,  
601 Ezequiel Gleichgerricht, Rutvik Harshad Desai, Nicholas Riccardi, Alexandra Basilakos, et al. Neural network  
602 bases of thematic semantic processing in language production. *Cortex*, 2022.
- 603 [11] Peter Brunner, Anthony L Ritaccio, Timothy M Lynch, Joseph F Emrich, J Adam Wilson, Justin C Williams,  
604 Erik J Aarnoutse, Nick F Ramsey, Eric C Leuthardt, Horst Bischof, et al. A practical procedure for real-time  
605 functional mapping of eloquent cortex using electrocorticographic signals in humans. *Epilepsy & Behavior*,  
606 15(3):278–286, 2009.
- 607 [12] Gyorgy Buzsaki. *Rhythms of the Brain*. Oxford University Press, 2006.
- 608 [13] György Buzsáki, Costas A Anastassiou, and Christof Koch. The origin of extracellular fields and currentseeg,  
609 ecog, lfp and spikes. *Nature reviews neuroscience*, 13(6):407–420, 2012.
- 610 [14] Ujwal Chaudhary, Niels Birbaumer, and Ander Ramos-Murguialday. Brain–computer interfaces for communica-  
611 tion and rehabilitation. *Nature Reviews Neurology*, 12(9):513, 2016.
- 612 [15] Jacquelyn A. Corley, Pouya Nazari, Vincent J. Rossi, Nora C. Kim, Louis F. Fogg, Thomas J. Hoeppepner, Travis R.  
613 Stoub, and Richard W. Byrne. Cortical stimulation parameters for functional mapping. *Seizure*, 45:36–41, 2017.
- 614 [16] Andrey Eliseyev and Tatiana Aksenova. Stable and artifact-resistant decoding of 3d hand trajectories from ecog  
615 signals using the generalized additive model. *Journal of neural engineering*, 11(6):066005, 2014.
- 616 [17] Michael J Fagan, Stephen R Ell, James M Gilbert, E Sarrazin, and Peter M Chapman. Development of a (silent)  
617 speech recognition system for patients following laryngectomy. *Medical engineering & physics*, 30(4):419–425,  
618 2008.
- 619 [18] Ananya Ganesh, Andre J Cervantes, and Philip R Kennedy. Slow firing single units are essential for optimal  
620 decoding of silent speech. *Frontiers in human neuroscience*, 16, 2022.
- 621 [19] Joris Guérin, Stephane Thiery, Eric Nyiri, Olivier Gibaru, and Byron Boots. Combining pretrained cnn feature  
622 extractors to enhance clustering of complex natural images. *Neurocomputing*, 423:551–571, 2021.

- 623 [20] Nicholas G Hatsopoulos and John P Donoghue. The science of neural interface systems. *Annual review of*  
624 *neuroscience*, 32:249–266, 2009.
- 625 [21] Stefan Haufe, Frank Meinecke, Kai Gorgen, Sven Dahne, John-Dylan Haynes, Benjamin Blankertz, and Felix  
626 Biemann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage*,  
627 87:96–110, 2014.
- 628 [22] Christian Herff, Lorenz Diener, Miguel Angrick, Emily Mugler, Matthew C Tate, Matthew A Goldrick, Dean J  
629 Krusienski, Marc W Slutzky, and Tanja Schultz. Generating natural, intelligible speech from brain activity in  
630 motor, premotor, and inferior frontal cortices. *Frontiers in neuroscience*, 13:1267, 2019.
- 631 [23] Christian Herff, Dean J Krusienski, and Pieter Kubben. The potential of stereotactic-*eeg* for brain-computer  
632 interfaces: current progress and future directions. *Frontiers in neuroscience*, 14:123, 2020.
- 633 [24] Ingo Hertrich, Susanne Dietrich, and Hermann Ackermann. The margins of the language network in the brain.  
634 *Frontiers in Communication*, 5:519955, 2020.
- 635 [25] Sepp Hochreiter and Jurgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,  
636 1997.
- 637 [26] Mark L Homer, Arto V Nurmikko, John P Donoghue, and Leigh R Hochberg. Sensors and decoding for intracor-  
638 tical brain computer interfaces. *Annual review of biomedical engineering*, 15:383–405, 2013.
- 639 [27] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional  
640 networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708,  
641 2017.
- 642 [28] Xuedong Huang, Alex Acero, Hsiao-Wuen Hon, and Raj Reddy. *Spoken Language Processing: A Guide to*  
643 *Theory, Algorithm, and System Development*. Prentice Hall PTR, USA, 1st edition, 2001.
- 644 [29] Prasanna Jayakar, Jean Gotman, A. Simon Harvey, Andre Palmi, Laura Tassi, Donald Schomer, Francois  
645 Dubeau, Fabrice Bartolomei, Alice Yu, Pavel Krek, Demetrios Velis, and Philippe Kahane. Diagnostic utility of  
646 invasive *eeg* for epilepsy surgery: Indications, modalities, and techniques. *Epilepsia*, 57(11):1735–1747, 2016.
- 647 [30] Rachel Kaye, Christopher G Tang, and Catherine F Sinclair. The electrolarynx: voice restoration after total  
648 laryngectomy. *Medical Devices (Auckland, NZ)*, 10:133, 2017.
- 649 [31] Christopher K Kovach, Naotsugu Tsuchiya, Hiroto Kawasaki, Hiroyuki Oya, Mathew A Howard III, and Ralph  
650 Adolphs. Manifestation of ocular-muscle *emg* contamination in human intracranial recordings. *Neuroimage*,  
651 54(1):213–233, 2011.
- 652 [32] Jan Kubanek, Peter Brunner, Aysegul Gunduz, David Poeppel, and Gerwin Schalk. The tracking of speech  
653 envelope in the human cortex. *PLoS one*, 8(1):e53398, 2013.
- 654 [33] Mikhail A Lebedev and Miguel AL Nicolelis. Brain-machine interfaces: From basic science to neuroprostheses  
655 and neurorehabilitation. *Physiological reviews*, 97(2):767–837, 2017.
- 656 [34] Sergio Machado, Fernanda Araujo, Flavia Paes, Bruna Velasques, Mario Cunha, Henning Budde, Luis F Basile,  
657 Renato Anghinah, Oscar Arias-Carrion, Mauricio Cagy, et al. *Eeg*-based brain-computer interfaces: an overview  
658 of basic concepts and clinical applications in neurorehabilitation. *Reviews in the Neurosciences*, 21(6):451–468,  
659 2010.
- 660 [35] Joseph N Mak and Jonathan R Wolpaw. Clinical applications of brain-computer interfaces: current state and  
661 future prospects. *IEEE reviews in biomedical engineering*, 2:187–199, 2009.
- 662 [36] Joseph G Makin, David A Moses, and Edward F Chang. Machine translation of cortical activity to text with an  
663 encoder–decoder framework. *Nature Neuroscience*, 23(4):575–582, 2020.
- 664 [37] L. Marple. A new autoregressive spectrum analysis algorithm. *IEEE Transactions on Acoustics, Speech, and*  
665 *Signal Processing*, 28:441–454, 1980.
- 666 [38] Brian McFee, Colin Raffel, Dawen Liang, Daniel P Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto.  
667 *librosa*: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*,  
668 volume 8, pages 18–25. Citeseer, 2015.

- 669 [39] David A Moses, Sean L Metzger, Jessie R Liu, Gopala K Anumanchipalli, Joseph G Makin, Pengfei F Sun, Josh  
670 Chartier, Maximilian E Dougherty, Patricia M Liu, Gary M Abrams, et al. Neuroprosthesis for decoding speech  
671 in a paralyzed person with anarthria. *New England Journal of Medicine*, 385(3):217–227, 2021.
- 672 [40] Emily M Mugler, James L Patton, Robert D Flint, Zachary A Wright, Stephan U Schuele, Joshua Rosenow,  
673 Jerry J Shih, Dean J Krusienski, and Marc W Slutzky. Direct classification of all american english phonemes  
674 using signals from functional speech motor cortex. *Journal of neural engineering*, 11(3):035015, 2014.
- 675 [41] Klaus-Robert Müller, Matthias Krauledat, Guido Dornhege, Gabriel Curio, and Benjamin Blankertz. Machine  
676 learning techniques for brain-computer interfaces. *Biomed. Tech*, 49(1):11–22, 2004.
- 677 [42] Luis Fernando Nicolas-Alonso and Jaime Gomez-Gil. Brain computer interfaces, a review. *Sensors*, 12(2):1211–  
678 1279, 2012.
- 679 [43] Thomas J Oxley, Nicholas L Opie, Sam E John, Gil S Rind, Stephen M Ronayne, Tracey L Wheeler, Jack W Judy,  
680 Alan J McDonald, Anthony Dornom, Timothy JH Lovell, et al. Minimally invasive endovascular stent-electrode  
681 array for high-fidelity, chronic recordings of cortical neural activity. *Nature biotechnology*, 34(3):320–327, 2016.
- 682 [44] Miguel Pais-Vieira, Mikhail Lebedev, Carolina Kunicki, Jing Wang, and Miguel Nicolelis. A brain-to-brain  
683 interface for real-time sharing of sensorimotor information. *Scientific reports*, 3:1319, 02 2013.
- 684 [45] Artur Petrosyan, Mikhail Sinkin, Mikhail Lebedev, and Alexei Ossadtchi. Decoding and interpreting cortical  
685 signals with a compact convolutional neural network. *Journal of Neural Engineering*, 18(2):026019, 2021.
- 686 [46] Nick F Ramsey, Efraim Salari, Erik J Aarnoutse, Mariska J Vansteensel, Martin G Bleichner, and ZV Freuden-  
687 burg. Decoding spoken phonemes from sensorimotor cortex with high-density ecog grids. *Neuroimage*, 180:301–  
688 311, 2018.
- 689 [47] Philémon Roussel, Gaël Le Godais, Florent Bocquelet, Marie Palma, Jiang Hongjie, Shaomin Zhang, Anne-  
690 Lise Giraud, Pierre Mégevand, Kai Miller, Johannes Gehrig, et al. Observation and assessment of acoustic  
691 contamination of electrophysiological brain signals during speech production and sound perception. *Journal of*  
692 *Neural Engineering*, 17(5):056028, 2020.
- 693 [48] Philémon Roussel, Gaël Le Godais, Florent Bocquelet, Marie Palma, Jiang Hongjie, Shaomin Zhang, Philippe  
694 Kahane, Stéphan Chabardès, and Blaise Yvert. Acoustic contamination of electrophysiological brain signals  
695 during speech production and sound perception. *BioRxiv*, page 722207, 2019.
- 696 [49] Gerwin Schalk and Eric C Leuthardt. Brain-computer interfaces using electrocorticographic signals. *IEEE*  
697 *reviews in biomedical engineering*, 4:140–154, 2011.
- 698 [50] Lauren J Silbert, Christopher J Honey, Erez Simony, David Poeppel, and Uri Hasson. Coupled neural systems  
699 underlie the production and comprehension of naturalistic narrative speech. *Proceedings of the National Academy*  
700 *of Sciences*, 111(43):E4687–E4696, 2014.
- 701 [51] M Sinkin, A Osadchiy, M Lebedev, K Volkova, M Kondratova, I Trifonov, et al. High resolution passive speech  
702 mapping in dominant hemisphere glioma surgery. *Russ. J. Neurosurg*, 21:12–18, 2019.
- 703 [52] Frank Soong and B Juang. Line spectrum pair (lsp) and speech data compression. In *ICASSP’84. IEEE Interna-*  
704 *tional Conference on Acoustics, Speech, and Signal Processing*, volume 9, pages 37–40. IEEE, 1984.
- 705 [53] Galina Spitsyna, Jane E Warren, Sophie K Scott, Federico E Turkheimer, and Richard JS Wise. Converging  
706 language streams in the human temporal lobe. *Journal of Neuroscience*, 26(28):7328–7336, 2006.
- 707 [54] Stanley Smith Stevens, John Volkman, and Edwin Broomell Newman. A scale for the measurement of the  
708 psychological magnitude pitch. *The journal of the acoustical society of america*, 8(3):185–190, 1937.
- 709 [55] Pengfei Sun, Gopala K Anumanchipalli, and Edward F Chang. Brain2char: a deep architecture for decoding text  
710 from brain recordings. *Journal of Neural Engineering*, 17(6):066015, 2020.
- 711 [56] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan,  
712 Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE*  
713 *conference on computer vision and pattern recognition*, pages 1–9, 2015.
- 714 [57] Ksenia Volkova, Mikhail A Lebedev, Alexander Kaplan, and Alexei Ossadtchi. Decoding movement from elec-  
715 trocorticographic activity: A review. *Frontiers in neuroinformatics*, 13:74, 2019.

- 716 [58] Sarah K Wandelt, Spencer Kellis, David A Bjånes, Kelsie Pejsa, Brian Lee, Charles Liu, and Richard A Andersen.  
717 Decoding grasp and speech signals from the cortical grasp circuit in a tetraplegic human. *Neuron*, 2022.
- 718 [59] Francis R Willett, Donald T Avansino, Leigh R Hochberg, Jaimie M Henderson, and Krishna V Shenoy. High-  
719 performance brain-to-text communication via handwriting. *Nature*, 593(7858):249–254, 2021.
- 720 [60] Guy H Wilson, Sergey D Stavisky, Francis R Willett, Donald T Avansino, Jessica N Kelemen, Leigh R Hochberg,  
721 Jaimie M Henderson, Shaul Druckmann, and Krishna V Shenoy. Decoding spoken english from intracortical  
722 electrode arrays in dorsal precentral gyrus. *Journal of Neural Engineering*, 17(6):066007, 2020.
- 723 [61] Min Xu, Ling-Yu Duan, Jianfei Cai, Liang-Tien Chia, Changsheng Xu, and Qi Tian. Hmm-based audio keyword  
724 generation. In *Pacific-Rim Conference on Multimedia*, pages 566–574. Springer, 2004.