

REPORT DOCUMENTATION PAGE			Form Approved OMB NO. 0704-0188		
<p>The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA, 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number. PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.</p>					
1. REPORT DATE (DD-MM-YYYY)		2. REPORT TYPE		3. DATES COVERED (From - To)	
		New Reprint		-	
4. TITLE AND SUBTITLE			5a. CONTRACT NUMBER		
Brain-to-text: decoding spoken phrases from phone representations in the brain			W911NF-14-1-0440		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
			611102		
6. AUTHORS			5d. PROJECT NUMBER		
Dominic Heger, Adriana de Pestors, Christian Herff, Dominic Telaar, Peter Brunner, Gerwin Schalk, Tanja Schultz					
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAMES AND ADDRESSES			8. PERFORMING ORGANIZATION REPORT NUMBER		
Health Research, Inc. @ Wadsworth 150 Broadway, Suite 560  Menands, NY 12204 -2719					
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS (ES)			10. SPONSOR/MONITOR'S ACRONYM(S) ARO		
U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211			11. SPONSOR/MONITOR'S REPORT NUMBER(S) 62347-LS.1		
12. DISTRIBUTION AVAILABILITY STATEMENT					
Approved for public release; distribution is unlimited.					
13. SUPPLEMENTARY NOTES					
The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.					
14. ABSTRACT					
It has long been speculated whether communication between humans and machines based on natural speech related cortical activity is possible. Over the past decade, studies have suggested that it is feasible to recognize isolated aspects of speech from neural signals, such as auditory features, phones or one of a few isolated words. However, until now it remained an unsolved challenge to decode continuously spoken speech from the neural substrate associated with speech and language processing. Here, we show for the first time that continuously spoken speech can be decoded into the expressed words from intracranial electrocorticography (ECoG) recordings. Specifically:					
15. SUBJECT TERMS					
ECoG; automatic speech recognition; brain-computer interface; broadband gamma; electrocorticography; pattern recognition; speech decoding; speech production					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	15. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON
a. REPORT	b. ABSTRACT	c. THIS PAGE	UU		Gerwin Schalk
UU	UU	UU			19b. TELEPHONE NUMBER
					518-486-2559

## **Report Title**

Brain-to-text: decoding spoken phrases from phone representations in the brain

### **ABSTRACT**

It has long been speculated whether communication between humans and machines based on natural speech related cortical activity is possible. Over the past decade, studies have suggested that it is feasible to recognize isolated aspects of speech from neural signals, such as auditory features, phones or one of a few isolated words. However, until now it remained an unsolved challenge to decode continuously spoken speech from the neural substrate associated with speech and language processing. Here, we show for the first time that continuously spoken speech can be decoded into the expressed words from intracranial electrocorticographic (ECoG) recordings. Specifically, we implemented a system, which we call Brain-To-Text that models single phones, employs techniques from automatic speech recognition (ASR), and thereby transforms brain activity while speaking into the corresponding textual representation. Our results demonstrate that our system can achieve word error rates as low as 25% and phone error rates below 50%. Additionally, our approach contributes to the current understanding of the neural basis of continuous speech production by identifying those cortical regions that hold substantial information about individual phones. In conclusion, the Brain-To-Text system described in this paper represents an important step toward human-machine communication based on imagined speech.

---

**REPORT DOCUMENTATION PAGE (SF298)**  
**(Continuation Sheet)**

---

Continuation for Block 13

ARO Report Number 62347.1-LS

Brain-to-text: decoding spoken phrases from ph...

Block 13: Supplementary Note

© 2015 . Published in Frontiers in Neuroscience, Vol. Ed. 0 9, (0) (2015), (, (0). DoD Components reserve a royalty-free, nonexclusive and irrevocable right to reproduce, publish, or otherwise use the work for Federal purposes, and to authorize others to do so (DODGARS §32.36). The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.

Approved for public release; distribution is unlimited.

## Brain-to-text: Decoding spoken phrases from phone representations in the brain

Christian Herff, Dominic Heger, Adriana de\_Pesters, Dominic Telaar, Peter Brunner, Gerwin Schalk and Tanja Schultz

Journal Name:	Frontiers in Neuroengineering
ISSN:	1662-6443
Article type:	Original Research Article
Received on:	09 Apr 2015
Accepted on:	18 May 2015
Provisional PDF published on:	18 May 2015
Frontiers website link:	<a href="http://www.frontiersin.org">www.frontiersin.org</a>
Citation:	Herff C, Heger D, De_pesters A, Telaar D, Brunner P, Schalk G and Schultz T(2015) Brain-to-text: Decoding spoken phrases from phone representations in the brain. <i>Front. Neuroeng.</i> 8:6. doi:10.3389/fneng.2015.00006
Copyright statement:	© 2015 Herff, Heger, De_pesters, Telaar, Brunner, Schalk and Schultz. This is an open-access article distributed under the terms of the <a href="http://creativecommons.org/licenses/by/2.0/">Creative Commons Attribution License (CC BY)</a> . The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after rigorous peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.



1

# Brain-to-text: Decoding spoken phrases from phone representations in the brain

Christian Herff<sup>1\*</sup>, Dominic Heger<sup>1\*</sup>, Adriana de Pestors<sup>2,4</sup>, Dominic Telaar<sup>1</sup>, Peter Brunner<sup>2,3</sup>, Gerwin Schalk<sup>2,3,4</sup>, Tanja Schultz<sup>1</sup>

<sup>1</sup> Cognitive Systems Lab, Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Karlsruhe, Germany

<sup>2</sup> National Center for Adaptive Neurotechnologies, Wadsworth Center, New York State Department of Health, Albany, NY, USA

<sup>3</sup> Department of Neurology, Albany Medical College, Albany, USA

<sup>4</sup> Department of Biomedical Sciences, State University of New York at Albany, Albany, NY, USA

Correspondence\*:

Christian Herff  
Cognitive Systems Lab, Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Adenauerring 4, 76131 Karlsruhe, Germany, christian.herff@kit.edu

Dominic Heger  
Cognitive Systems Lab, Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Adenauerring 4, 76131 Karlsruhe, Germany, dominic.heger@kit.edu

\* These authors contributed equally to this work.

## 2 ABSTRACT

3 It has long been speculated whether communication between humans and machines based on  
4 natural speech related cortical activity is possible. Over the past decade, studies have suggested  
5 that it is feasible to recognize isolated aspects of speech from neural signals, such as auditory  
6 features, phones or one of a few isolated words. However, until now it remained an unsolved  
7 challenge to decode continuously spoken speech from the neural substrate associated with  
8 speech and language processing. Here, we show for the first time that continuously spoken  
9 speech can be decoded into the expressed words from intracranial electrocorticographic (ECoG)  
10 recordings. Specifically, we implemented a system, which we call *Brain-To-Text* that models  
11 single phones, employs techniques from automatic speech recognition (ASR), and thereby  
12 transforms brain activity while speaking into the corresponding textual representation. Our  
13 results demonstrate that our system achieved word error rates as low as 25% and phone  
14 error rates below 50%. Additionally, our approach contributes to the current understanding  
15 of the neural basis of continuous speech production by identifying those cortical regions that  
16 hold substantial information about individual phones. In conclusion, the Brain-To-Text system  
17 described in this paper represents an important step towards human-machine communication  
18 based on imagined speech.

19 **Keywords:** electrocorticography, ECoG, speech production, automatic speech recognition, brain-computer interface, speech  
20 decoding, pattern recognition, broadband gamma

## 1 INTRODUCTION

21 Communication with computers or humans by thought alone, is a fascinating concept and has long been  
22 a goal of the brain-computer interface (BCI) community (**Wolpaw et al. (2002)**). Traditional BCIs use  
23 motor imagery (**McFarland et al. (2000)**) to control a cursor or to choose between a selected number of  
24 options. Others use event-related potentials (ERPs) (**Farwell and Donchin (1988)**) or steady-state evoked  
25 potentials (**Sutter (1992)**) to spell out texts. These interfaces have made remarkable progress in the last  
26 years, but are still relatively slow and unintuitive. The possibility of using covert speech, i.e. imagined  
27 continuous speech processes recorded from the brain for human-computer communication may improve  
28 BCI communication speed and also increase their usability. Numerous members of the scientific  
29 community, including linguists, speech processing technologists, and computational neuroscientists have  
30 studied the basic principles of speech and analyzed its fundamental building blocks. However, the high  
31 complexity and agile dynamics in the brain make it challenging to investigate speech production with  
32 traditional neuroimaging techniques. Thus, previous work has mostly focused on isolated aspects of  
33 speech in the brain.

34 Several recent studies have begun to take advantage of the high spatial resolution, high temporal  
35 resolution and high signal-to-noise ratio of signals recorded directly from the brain (electrocorticography  
36 (ECoG)). Several studies used ECoG to investigate the temporal and spatial dynamics of speech perception  
37 (**Kubanek et al. (2013)**; **Canolty et al. (2007)**). Other studies highlighted the differences between  
38 receptive and expressive speech areas (**Towle et al. (2008)**; **Fukuda et al. (2010)**). Further insights into  
39 the isolated repetition of phones and words has been provided in (**Leuthardt et al. (2011b)**; **Pei et al.**  
40 **(2011b)**). **Pasley et al. (2012)** showed that auditory features of perceived speech could be reconstructed  
41 from brain signals. In a study with a completely paralyzed subject, **Guenther et al. (2009)** showed that  
42 brain signals from speech-related regions could be used to synthesize vowel formants. Following up on  
43 these results, **Martin et al. (2014)** decoded spectrotemporal features of overt and covert speech from  
44 ECoG recordings. Evidence for a neural representation of phones and phonetic features during speech  
45 perception was provided in **Chang et al. (2010)** and **Mesgarani et al. (2014)**, but these studies did not  
46 investigate continuous speech production. Other studies investigated the dynamics of the general speech  
47 production process (**Crone et al. (2001a,b)**). A large number of studies have classified isolated aspects  
48 of speech processes for communication with or control of computers. **Deng et al. (2010)** decoded three  
49 different rhythms of imagined syllables. Neural activity during the production of isolated phones was  
50 used to control a one-dimensional cursor accurately (**Leuthardt et al. (2011a)**). **Formisano et al. (2008)**  
51 decoded isolated phones using functional magnetic resonance imaging (fMRI). Vowels and consonants  
52 were successfully discriminated in limited pairings in **Pei et al. (2011a)**. **Blakely et al. (2008)** showed  
53 robust classification of four different phonemes. Other ECoG studies classified syllables (**Bouchard**  
54 **and Chang (2014)**) or a limited set of words (**Kellis et al. (2010)**). Extending this idea, the imagined  
55 production of isolated phones was classified in **Brumberg et al. (2011)**. Recently, **Mugler et al. (2014b)**  
56 demonstrated the classification of a full set of phones within manually segmented boundaries during  
57 isolated word production.

58 To make use of these promising results for BCIs based on continuous speech processes, the analysis  
59 and decoding of isolated aspects of speech production has to be extended to continuous and fluent speech  
60 processes. While relying on isolated phones or words for communication with interfaces would improve  
61 current BCIs drastically, communication would still not be as natural and intuitive as continuous speech.  
62 Furthermore, to process the content of the spoken phrases, a textual representation has to be extracted  
63 instead of a reconstruction of acoustic features. In our present study, we address these issues by analyzing  
64 and decoding brain signals during continuously produced overt speech. This enables us to reconstruct  
65 continuous speech into a sequence of words in textual form, which is a necessary step towards human-  
66 computer communication using the full repertoire of imagined speech. We refer to our procedure that  
67 implements this process as *Brain-to-Text*. *Brain-to-Text* implements and combines understanding from  
68 neuroscience and neurophysiology (suggesting the locations and brain signal features that should be  
69 utilized), linguistics (phone and language model concepts), and statistical signal processing and machine  
70 learning. Our results suggest that the brain encodes a repertoire of phonetic representations that can be

71 decoded continuously during speech production. At the same time, the neural pathways represented within  
72 our model offer a glimpse into the complex dynamics of the brain's fundamental building blocks during  
73 speech production.

## 2 MATERIAL & METHODS

### 2.1 SUBJECTS

74 Seven epileptic patients at Albany Medical Center (Albany, New York, USA) participated in this study.  
75 All subjects gave informed consent to participate in the study, which was approved by the Institutional  
76 Review Board of Albany Medical College and the Human Research Protections Office of the US Army  
77 Medical Research and Materiel Command. Relevant patient information is given in Figure 1.

### 2.2 ELECTRODE PLACEMENT

78 Electrode placement was solely based on clinical needs of the patients. All subjects had electrodes  
79 implanted on the left hemisphere and covered relevant areas of the frontal and temporal lobes. Electrode  
80 grids (Ad-Tech Medical Corp., Racine, WI; PMT Corporation, Chanhassen, MN) were composed of  
81 platinum-iridium electrodes (4 mm in diameter, 2.3 mm exposed) embedded in silicon with an inter-  
82 electrode distance of 0.6-1 cm. Electrode positions were registered in a post-operative CT scan and  
83 co-registered with a pre-operative MRI scan. Figure 1 shows electrode positions of all 7 subjects and the  
84 combined electrode positions. To compare average activation patterns across subjects, we co-registered  
85 all electrode positions in common Talairach space. We rendered activation maps using the NeuralAct  
86 software package (**Kubaneck and Schalk (2014)**).

87 **Figure 1.** Electrode positions for all seven subjects. Captions include age (years old (y/o)) and sex  
88 of subjects. Electrode locations were identified in a post-operative CT and co-registered to preoperative  
89 MRI. Electrodes for subject 3 are on an average Talairach brain. Combined electrode placement in joint  
90 Talairach space for comparison of all subjects. Participant 1 (yellow), subject 2 (magenta), subject 3  
91 (cyan), subject 5 (red), subject 6 (green) and subject 7 (blue). Participant 4 was excluded from joint  
92 analysis as the data did not yield sufficient activations related to speech activity (see Section 2.4).

### 2.3 EXPERIMENT

93 We recorded brain activity during speech production of seven subjects using electrocorticographic (ECoG)  
94 grids that had been implanted as part of presurgical procedures preparatory to epilepsy surgery. ECoG  
95 provides electrical potentials measured directly on the brain surface at a high spatial and temporal  
96 resolution, unfiltered by skull and scalp. ECoG signals were recorded by BCI2000 (**Schalk et al. (2004)**)  
97 using eight 16-channel g.USBamp biosignal amplifiers (g.tec, Graz, Austria). In addition to the electrical  
98 brain activity measurements, we recorded the acoustic waveform of the subjects' speech. Participant's  
99 voice data was recorded with a dynamic microphone (Samson R21s) and digitized using a dedicated  
100 g.USBamp in sync with the ECoG signals. The ECoG and acoustic signals were digitized at a sampling  
101 rate of 9600 Hz.

102 During the experiment, text excerpts from historical political speeches (i.e., Gettysburg Address (**Roy**  
103 **and Basler (1955)**), JFK's Inaugural Address (**Kennedy (1989)**), a childrens' story (**Crane et al. (1867)**)  
104 or *Charmed* fan-fiction (**unknown (2009)**) were displayed on a screen in about 1 m distance from the  
105 subject. The texts scrolled across the screen from right to left at a constant rate. This rate was adjusted  
106 to be comfortable for the subject prior to the recordings (rate of scrolling text: 42-76 words/min). During  
107 this procedure, subjects were familiarized with the task.

108 Each subject was instructed to read the text aloud as it appeared on the screen. A session was repeated  
109 2 – 3 times depending on the mental and physical condition of the subjects. Table 1 summarizes data

110 recording details for every session. Since the amount of data of the individual sessions of subject 2 is very  
 111 small, we combined all three sessions of this subject in the analysis.

**Table 1.** Data recording details for every session.

Participant	Session	Text	Number of phrases	Total recording length (s)
1	1	Gettysburg address	36	279.87
	2	JFK inaugural	38	326.90
2	1	Humpty Dumpty	21	129.87
	2	Humpty Dumpty	21	129.07
	3	Humpty Dumpty	21	126.37
3	1	Charmed fan-fiction	42	310.27
	2	Charmed fan-fiction	40	310.93
	3	Charmed fan-fiction	41	307.50
4	1	Gettysburg address	38	299.67
	2	Gettysburg address	38	311.97
5	1	JFK inaugural	49	341.77
	2	Gettysburg address	39	222.57
6	1	Gettysburg address	38	302.83
7	1	JFK inaugural	48	590.10
	2	Gettysburg address	38	391.43

112 We cut the read-out texts of all subjects into 21 to 49 phrases, depending on the session length, along  
 113 pauses in the audio recording. The audio recordings were phone-labeled using our in-house speech  
 114 recognition toolkit BioKIT (Telaar et al. (2014)) (see Section 2.5). Because the audio and ECoG data  
 115 were recorded in synchronization (see Figure 2), this procedure allowed us to identify the ECoG signals  
 116 that were produced at the time of any given phones. Figure 2 shows the experimental setup and the phone  
 117 labeling.

118 **Figure 2.** Synchronized recording of ECoG and acoustic data. Acoustic data are labeled using our in-  
 119 house decoder BioKIT, i.e. the acoustic data samples are assigned to corresponding phones. These phone  
 120 labels are then imposed on the neural data.

## 2.4 DATA PRE-SELECTION

121 In an initial data pre-selection, we tested whether speech activity segments could be distinguished from  
 122 those with no speech activity in ECoG data. For this purpose, we fitted a multivariate normal distribution  
 123 to all feature vectors (see Section 2.6 for a description of the feature extraction) containing speech activity  
 124 derived from the acoustic data and one to feature vectors when the subject was not speaking. We then  
 125 determined whether these models could be used to classify general speech activity above chance level,  
 126 applying a leave-one-phrase-out validation.

127 Based on this analysis, both sessions of subject 4 and session 2 of subject 5 were rejected, as they did not  
 128 show speech related activations that could be classified significantly better than chance (t-test,  $p > 0.05$ ).  
 129 To compare against random activations without speech production, we employed the same randomization  
 130 approach as described in Section 2.11.

## 2.5 PHONE LABELING

131 Phone labels of the acoustic recordings were created in a three-step process using an English automatic  
 132 speech recognition (ASR) system trained on broadcast news. First, we calculated a Viterbi forced  
 133 alignment (Huang et al. (2001)), which is the most likely sequence of phones for the acoustic data samples  
 134 given the words in the transcribed text and the acoustic models of the ASR system. In a second step, we  
 135 adapted the Gaussian mixture model (GMM)-based acoustic models using maximum likelihood linear  
 136 regression (MLLR) (Gales (1998)). This adaptation was performed separately for each session to obtain  
 137 session-dependent acoustic models specialized to the signal and speaker characteristics, which is known  
 138 to increase ASR performance. We estimated a MLLR transformation from the phone sequence computed  
 139 in step one and used only those segments which had a high confidence score that the segment was emitted  
 140 by the model attributed to them. Third, we repeated the Viterbi forced alignment using each session's  
 141 adapted acoustic models yielding the final phone alignments. The phone labels calculated on the acoustic  
 142 data are then imposed on the ECoG data.

143 Due to the very limited amount of training data for the neural models, we reduced the amount of distinct  
 144 phone types and grouped similar phones together for the ECoG models. The grouping was based on  
 145 phonetic features of the phones. See Table 2 for the grouping of phones.

**Table 2.** Grouping of phones. English phones are based on the International Phonetic Alphabet (IPA).

Grouped Phone	IPA phones
aa	ɑ æʌ
b	b
ch	tʃ ʃ ʒ
eh	ɛ ɜː eɪ
f	f
hh	h
ih	i ɪ
jh	dʒ ʒ j
k	k
l	l
m	m
n	n ŋ
ow	oʊ ɔ
p	p
r	r
s	s z ʒ θ
t	t d
uw	u ʊ
v	v
w	w
Diphthongs	
ow ih	ɔɪ
aa ih	aɪ
aa ow	aʊ

## 2.6 FEATURE EXTRACTION

146 We segmented the neural signal data continuously into 50 ms intervals with an overlap of 25 ms, which  
 147 enabled us to capture the fast cortical processes underlying phones, while being long enough to extract

148 broadband (70–170 Hz) gamma activity reliably. Each of the 50 ms intervals was labeled with the  
 149 corresponding phone obtained from the audio phone labeling. We extracted broadband-gamma activations  
 150 as they are known to be highly task-related for motor tasks (Miller et al. (2007)), music perception (Potes  
 151 et al. (2012)), auditory processes (Pasley et al. (2012); Pei et al. (2011b)) and word repetition (Leuthardt  
 152 et al. (2011b)). Broadband-gamma features were extracted from the ECoG electrical potentials as follows:  
 153 linear trends in the raw signals were removed from each channel. The signals were down-sampled from  
 154 9600 Hz to 600 Hz sampling rate. Channels strongly affected by noise were identified and excluded from  
 155 further processing. Specifically, we calculated the energy in the frequency band 58-62 Hz (line noise) and  
 156 removed channels with more noise energy than two interquartile ranges above the third quartile of the  
 157 energy of all channels in the data set. This way, an average of 7.0 (std 6.5) channels were removed per  
 158 subject.

159 The remaining channels were re-referenced to a common average (i.e., CAR filtering). Elliptic IIR low-  
 160 pass and high-pass filters were applied to represent broadband gamma activity in the signals. An elliptic  
 161 IIR notch filter (118-122 Hz, filter order 13) was applied to attenuate the first harmonic of 60 Hz line  
 162 noise, which is within the broadband gamma frequency range.

163 Resulting 50 ms intervals are denoted as  $X_{i,c}(t)$  and consist of  $n$  samples ( $t \in [1, \dots, n]$ ). For each  
 164 interval  $i$  and channel  $c$ , the signal energy  $E_{i,c}$  was calculated and the logarithm was applied to make the  
 165 distribution of the energy features approximately Gaussian:  $E_{i,c} = \log(\frac{1}{n} \sum_{t=1}^n X_{i,c}(t)^2)$ . The logarithmic  
 166 broadband gamma power of all channels were concatenated into one feature vector  $E_i = [E_{i,1}, \dots, E_{i,d}]$ .  
 167 To integrate context information and temporal dynamics of the neural activity for each interval, we  
 168 included neighboring intervals up to 200 ms prior to and after the current interval, similar context sizes  
 169 have been found relevant in speech perception studies (Sahin et al. (2009)). Therefore, each feature vector  
 170 was stacked with four feature vectors in the past and four feature vectors in the future. Stacked feature  
 171 vectors  $F_i = [E_{i-4}, \dots, E_i, \dots, E_{i+4}]^T$  were extracted every 25 ms over the course of the recording  
 172 sessions and the fitting phone label (ground truth from acoustic phone labeling) was associated.

## 2.7 IDENTIFICATION OF DISCRIMINABILITY

173 The high temporal and spatial resolution of ECoG recordings allowed us to trace the temporal dynamics  
 174 of speech production through the areas in the brain relevant for continuous natural speech production. To  
 175 investigate such cortical regions of high relevance, we calculated the mean symmetrized Kullback-Leibler  
 176 divergence (KL-div) among the phone models for each electrode position and at every time interval.

177 The Kullback-Leibler divergence (KL-div) is a measure of the difference between two distributions  $P$   
 178 and  $Q$ . It can be interpreted as the amount of discriminability between the neural activity models in bits. It  
 179 is non-symmetric and does not satisfy the triangle inequality. The KL-div can be interpreted as the amount  
 180 of extra bits needed to code samples from  $P$  when using  $Q$  to estimate  $P$ . When both distributions  $P$  and  
 181  $Q$  are normal distributions with means  $\mu_0$  and  $\mu_1$  and covariances  $\Sigma_0$  and  $\Sigma_1$ , respectively, the KL-div can  
 182 be easily calculated as

$$D_{KL}(N_0||N_1) = \frac{1}{2}(tr(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0) - d - \log_2(\frac{\det(\Sigma_0)}{\det(\Sigma_1)})) \quad (1)$$

183 with  $d$  being the dimensionality of the distributions. The closed-form of the KL-div enables us to calculate  
 184 the difference between two phone models. To estimate the discriminability of a feature  $E_{i,c}$  (log broadband  
 185 gamma power of a particular channel and time interval) for the classification of phones, we calculate the  
 186 mean KL-div between all pairs of phones for this particular feature. The mean between all divergences  
 187 symmetrizes the KL-div and yields one number in bits as the estimation of the discriminability of this  
 188 particular feature  $E_{i,c}$ .

## 2.8 FEATURE SELECTION

189 We selected features with the largest average distance between phone models based on the mean KL-  
 190 div (cf. previous section) in the training data during each run of the leave-one-phrase-out validation.  
 191 The number of features selected was automatically determined based on the distribution of KL-div for  
 192 this specific run as follows: We normalized the mean KL-div values  $d_k$  for every feature  $k$  by their  
 193 average ( $\hat{d}_k = \frac{d_k}{\sum_k d_k}$ ). Then, we sorted the values in descending order and selected features with large  
 194 normalized mean KL-div until the sorted sequence did not decline more than a threshold  $t = -0.05$ :  
 195  $\arg \max_l \text{sort}(\hat{d}_k)_l - \text{sort}(\hat{d}_k)_{l+1} < t$ . The threshold value  $t = -0.05$  corresponds to a very low decline  
 196 in KL-div and thus reflected the point after which little additional information was present. This way, only  
 197 the  $l$  most relevant features are selected to limit the feature space.

198 Note that features are selected solely based on the Kullback-Leibler divergence in the training data and  
 199 do not include any prior assumptions on the suitability of specific regions for phone discrimination. We  
 200 further reduced the feature space dimensionality by linear discriminant analysis (LDA) (**Haeb-Umbach**  
 201 **and Ney** (1992)) using the phone labels on the training data.

## 2.9 ECOG PHONE MODEL TRAINING

202 Each phone was modeled in the extracted feature space by a normal distribution. Thus, models  
 203 characterized the mean contribution and variance of the neural activity measured at each electrode. We  
 204 represented the stacked cortical activity feature vectors  $F_i$  of each phone  $j$  by a model  $\lambda_j$  as a multivariate  
 205 Gaussian probability density function  $p(F_i|\lambda_j) \sim \mathcal{N}(\mu_j, \Sigma_j)$  determined by the mean feature vectors  $\mu_j$   
 206 and their diagonal variance matrix  $\Sigma_j$  calculated from training data. Gaussian models were chosen as they  
 207 represent the underlying feature distribution suitably well. Furthermore, Gaussian models can be robustly  
 208 calculated from a small amount of data, they are computationally very efficient and allow a closed form  
 209 calculation of the Kullback-Leibler-Divergence.

## 2.10 DECODING APPROACH

210 Following a common idea of modern speech recognition technology (**Rabiner** (1989); **Schultz and**  
 211 **Kirchhoff** (2006)), we combined the information about the observed neural activity with statistical  
 212 language information during the decoding process by Bayesian updating (**Rabiner** (1989)). Simplified,  
 213 the process can be understood (**Gales and Young** (2008)) as finding the sequence of words  $W = w_1 \dots w_L$   
 214 which is most likely given the observed ECoG feature segments  $X = F_1 \dots F_T$ . This probability  $P(W|X)$   
 215 can be transformed using Bayes' rule:

$$\hat{W} = \arg \max_W \{P(W|X)\} = \arg \max_W \{p(X|W)P(W)\} \quad (2)$$

216 Here, the likelihood  $p(X|W)$  is given by the ECoG phone models and  $P(W)$  is calculated using a  
 217 language model. The likelihood of ECoG phone models  $p(X|W)$  given a word  $W$  is calculated by  
 218 concatenating ECoG phone models to form words as defined in a pronunciation dictionary. Specifically,  
 219 we employed a pronunciation dictionary containing the mapping of phone sequences to words, for  
 220 example, describing that the word 'liberty' comprises of the phone sequence '/l/ /ih/ /b/ /er/ /t/ /iy/'.  
 221 We constructed a minimized and determinized search graph consisting of the phone sequences for each  
 222 recognizable word. To capture important syntactic and semantic information of language, we used a  
 223 statistical language model (**Jelinek** (1997); **Stolcke** (2002)) that predicts the next word given the preceding  
 224 words. In N-gram language modeling, this is done by calculating probabilities of single words and  
 225 probabilities for predicting words given the  $n - 1$  previous words. Probabilities for single word occurrence  
 226 ( $n = 1$ ) are called uni-grams. Probabilities for the co-occurrence of two words ( $n = 2$ ) are called bi-grams.  
 227 For the *Brain-to-Text* system, we estimate bi-grams on the texts read by the subjects. It is important to

228 note that even though this results in very specialized models, the correctness of our results is still assured,  
229 as the same language models are utilized for both the real as well as for the control analyses.

230 Finally, the decoding of spoken phrases from neural data  $X$  is performed by finding the word sequence  
231  $\hat{W}$  in the search graph that has the highest likelihood for producing the neural data with respect to the  
232 ECoG phone models and language information given by pronunciation dictionary and language model.

233 Figure 3 illustrates the different steps of decoding continuously spoken phrases from neural data. *ECoG*  
234 *signals over time* are recorded at every electrode and divided into 50ms segments. For each 50 ms interval of  
235 recorded *broadband gamma activity*, stacked feature vectors are calculated (*Signal processing*). For  
236 each *ECoG phone model* calculated on the training data, the likelihood that this model emitted a segment  
237 of ECoG features can be calculated, resulting in *phone likelihoods over time*. Combining these Gaussian  
238 *ECoG phone models* with language information in the form of a *dictionary* and an *n-gram language*  
239 *model*, the *Viterbi* algorithm calculates the *most likely word sequence* and corresponding *phone sequence*.  
240 To visualize the decoding path, the *most likely phone sequence* can be shown in the *phone likelihoods over*  
241 *time* (red marked areas). The system outputs the decoded word sequence. Overall, the system produces a  
242 textual representation from the measured brain activity (see also Supplementary Video).

243 **Figure 3.** Overview of the *Brain-to-Text* system: ECoG broadband gamma activities (50ms segments)  
244 for every electrode are recorded. Stacked broadband gamma features are calculated (Signal processing).  
245 Phone likelihoods over time can be calculated by evaluating all Gaussian ECoG phone models for every  
246 segment of ECoG features. Using ECoG phone models, a dictionary and an n-gram language model,  
247 phrases are decoded using the Viterbi algorithm. The most likely word sequence and corresponding phone  
248 sequence are calculated and the phone likelihoods over time can be displayed. Red marked areas in the  
249 phone likelihoods show most likely phone path. See also Supplementary video.

## 2.11 EVALUATION

250 For the evaluation of our *Brain-to-Text* system, we trained neural phone models using all but one phrase  
251 of a recording session and decoded the remaining phrase. This evaluation process was repeated for  
252 each phrase in the session. Through this leave-one-phrase-out validation, we make sure that all feature  
253 selection, dimensionality reduction and training steps are only performed on the training data while the  
254 test data remains completely unseen. For comparison, we performed the decoding with randomized phone  
255 models. This is a baseline that quantifies how well the language model and dictionary decode phrases  
256 without any neural information. To obtain an estimate for chance levels in our approach, we shifted  
257 the training data by half its length in each iteration of the leave-one-phrase-out validation while the  
258 corresponding labels remained unchanged. This way, the data for the random comparison models still  
259 have the typical properties of ECoG broadband gamma activity, but do not correspond to the underlying  
260 labels. Furthermore, as the labels are not changed, prior probabilities remain the same for the random and  
261 the actual model case. As the shifting point is different for all iterations of the specific session, we get an  
262 estimate of the chance level performance for every phrase. The mean over all these results thus allows a  
263 robust estimation of the true chance level (randomization test).

264 It is also important to bear in mind that *Brain-to-Text* is still at a disadvantage compared to traditional  
265 speech recognition systems as our data contained only several minutes of ECoG signals for each subject.  
266 This limited model complexity compared to traditional speech recognition systems, which are usually  
267 trained on thousands of hours of acoustic data and billions of words for language model training.

268 We evaluated the performance of our *Brain-to-Text* system with different dictionary sizes. For this  
269 purpose, we created new dictionaries for every test phrase including the words that were actually spoken  
270 plus a set of randomized set of words from the full dictionary. Created dictionaries were the same for  
271 *Brain-To-Text* and randomized models to ensure that the words chosen had no influence on the comparison  
272 between models. The language model was limited to the words in the dictionary accordingly. This  
273 approach allowed us to perpetually increase the dictionary size.

### 3 RESULTS

#### 3.1 REGIONS OF DISCRIMINABILITY

274 Figure 4 illustrates the spatio-temporal dynamics of the mean KL-div between the phone models on a joint  
275 brain surface (Talairach model (Talairach and Tournoux (1988))) for nine temporal intervals with co-  
276 registered electrodes of all subjects. KL-div values plotted in Figure 4 exceed 99% of the KL-div values  
277 with a randomized phone-alignment (data shifted by half its length while the labels remain the same).

278 **Figure 4.** Mean Kullback-Leibler Divergences between models for every electrode position of every  
279 subject. Combined electrode montage of all subjects except subject 4 in common Talairach space. Heat  
280 maps on rendered average brain shows regions of high discriminability (red). All shown discriminability  
281 exceeds chance level (larger than 99% of randomized discriminabilities). The temporal course of regions  
282 with high discriminability between phone models shows early differences in diverse areas up to 200 ms  
283 before the actual phone production. Phone models show high discriminability in sensorimotor cortex 50  
284 ms before production and yield different models in auditory regions of the superior temporal gyrus 100  
285 ms after production.

286 Starting 200 ms before the actual phone production, we see high KL-div values in diverse areas  
287 including Broca's area, which is generally associated with speech planning (Sahin et al. (2009)).  
288 150 ms prior to the phone production, Broca's area still has high KL-div scores, but additionally  
289 sensorimotor areas and regions in the superior temporal gyrus associated with auditory and language  
290 function show increasing discriminability. Subsequently, activations in Broca's area vanish and motor  
291 area discriminability increases until peaking at the interval between 0 and 50 ms (which corresponds to  
292 the average length of phones). Discriminability increases in auditory regions until approximately 150 ms  
293 after phone production.

#### 3.2 DECODING RESULTS

294 For each phrase to be decoded, the most likely phone-path can be efficiently calculated using Viterbi  
295 decoding (Rabiner (1989)). Comparing the extracted phone labels for each feature vector with the  
296 baseline labels from the audio alignment, we calculate single-frame accuracies for the decoding of phones  
297 from continuous speech production. Reducing the size of the dictionary to 10 words, including those that  
298 are to be evaluated, *Brain-to-Text* yielded significantly higher accuracies (two-sided t-test,  $p < 0.05$  for  
299 all sessions) for single phone decoding in all sessions compared to random models. Figure 5 (a) shows  
300 average phone recognition accuracies (green) and average random recognition accuracies (orange) for  
301 each session. The best session resulted in average accuracies above 50% for the correct classification of 20  
302 phones plus SILENCE. While all sessions resulted in significantly higher accuracies than random models,  
303 the results of subject 2 and subject 7 clearly outperform those of all other subjects. The outstanding  
304 performance of subject 7 might be explained by the high-density grid on the superior temporal gyrus.  
305 We further investigate the results of subject 7, session 1 (results for all other subjects and sessions can  
306 be found in the Supplementary Material) by investigating the confusion matrix (Figure 5 (b)) that shows  
307 which phones in the reference corresponded to which phones in the predicted phrase. The clearly visible  
308 diagonal in this confusion matrix illustrates that our approach reliably decodes the complete set of phones.

309 In *Brain-to-Text*, we decode entire word sequences of each test phrase. Even with a small dictionary  
310 size, a large number of different phrases can be produced, as the number of words may vary and words  
311 can be arbitrarily combined. Therefore, we utilize the Word Error Rate (WER) to measure the quality of a  
312 decoded phrase. The word error rate (WER) between a predicted phrase and the corresponding reference  
313 phrase consists of the number of editing steps in terms of substitutions, deletions and insertions of words  
314 necessary to produce the predicted phrase from the reference, divided by the amount of words in the  
315 reference.

316 Figure 5 (c) shows the average WER depending on dictionary size (green line). For all dictionary sizes,  
317 the performance is significantly better than randomized results (orange line). Significance was analyzed

318 using paired t-tests between the Word Error Rates of *Brain-To-Text* and the randomized models ( $p <$   
319  $0.001$ , one-sided paired t-test). With 10 words in the dictionary, 75% of all words are recognized correctly.  
320 The approach scales well for increasing dictionary sizes. Average phone true positive rates remain rather  
321 stable even when dictionary sizes increase (bars in Figure 5 (c)).

322 **Figure 5.** Results: (A) Frame-wise accuracy for all sessions. All sessions of all subjects show  
323 significantly higher true positive rates for *Brain-To-Text* (green bars) than for the randomized models  
324 (orange bars). (B) Confusion matrix for subject 7, session 1. The clearly visible diagonal indicates that  
325 all phones are decoded reliably. (C) Word Error Rates depending on dictionary size (lines). Word error  
326 rates for *Brain-To-Text* (green line) are lower than the randomized models for all dictionary sizes. Average  
327 true-positive rates across phones depending on dictionary size (bars) for subject 7, session 1. Phone true  
328 positive rates remain relatively stable for all dictionary sizes and are always much higher for *Brain-To-Text*  
329 than for the randomized models.

## 4 DISCUSSION

### 4.1 ECOG PHONE MODELS

330 Gaussian models as a generative statistical representation for log-transformed broadband gamma power  
331 have been found well-suited for the observed cortical activity (e.g. **Gasser et al.** (1982); **Crone et al.**  
332 (2001b)). These models facilitate the analysis of the spatial and temporal characteristics of each  
333 phone model within its 450 ms context. Note that the modeling of phones does not contradict recent  
334 findings of articulatory features in neural recordings during speech perception (**Mesgarani et al.** (2014);  
335 **Pulvermüller et al.** (2006)) and production (**Bouchard et al.** (2013); **Lotte et al.** (2015)), since multiple  
336 representations of the same acoustic phenomenon are likely.

337 Note that only one context-independent model is trained for each phone, i.e., without consideration of  
338 preceding or succeeding phones due to the limited amount of data, even though effects of context have  
339 been shown in neural data (**Mugler et al.** (2014a)). While context dependent modeling is very common  
340 in acoustic speech recognition (**Lee** (1990)) and known to significantly improve recognition performance,  
341 it requires substantially more training data than available in our ECoG setting.

### 4.2 REGIONS OF DISCRIMINABILITY

342 In our approach, the phone representation through Gaussian models allows for detailed analysis of cortical  
343 regions, which have high discriminability among the different phones over time. The cortical locations  
344 identified using the KL-div criterion are in agreement with those that have been identified during speech  
345 production and perception in isolated phoneme or word experiments (**Leuthardt et al.** (2011a); **Canolty**  
346 **et al.** (2007)). These findings extend the state-of-the-art by showing for the first time the dynamics for  
347 single phone discriminability and decoding during continuous speech production.

348 As our experiments demand overt speech production from prompted texts, it is evident that multiple  
349 processes are present in the recorded neural data, including speech production, motor actions, auditory  
350 processing, and language understanding. By demonstrating that phones can be discriminated from each  
351 other, we show that such a phone-based representation is indeed a viable form of modeling cortical activity  
352 of continuous speech in this mixture of activation patterns.

### 4.3 DECODING RESULTS

353 The reported phone decoding accuracies are significantly higher for *Brain-to-Text* than for randomized  
354 models in all subjects, which shows that continuous speech production can be modeled based on phone  
355 representations. The clearly visible diagonal in the confusion matrix Figure 5 (B) emphasizes that the  
356 decoding performance is based on a reliable detection of all phones and not only on a selected subset.

357 Different conditions, such as varying task performance of the subjects, and different positions and  
358 densities of the electrode grids, yielded highly variable decoding performances for the different subjects,  
359 however the low WER (see Supplementary Material) and phone true positive rates for subject 1,2 and 7  
360 imply the potential of *Brain-to-Text* for brain-computer interfaces.

#### 4.4 CONCLUSION

361 Decoding overt speech production is a necessary first step towards human-computer interaction through  
362 imagined speech processes. Our results show that with a limited set of words in the dictionary, *Brain-*  
363 *to-Text* reconstructs spoken phrases from neural data. The computational phone models in combination  
364 with language information make it possible to reconstruct words in unseen spoken utterances solely based  
365 on neural signals (see Supplementary Video). Despite the fact that the evaluations in this article have  
366 been performed offline, all processing steps of *Brain-to-Text* and the decoding approach are well suited  
367 for eventual real-time online application on desktop computers. The approach introduced here may have  
368 important implications for the design of novel brain-computer interfaces, because it may eventually allow  
369 people to communicate solely based on brain signals associated with natural language function and with  
370 scalable vocabularies.

#### DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

371 The authors declare that the research was conducted in the absence of any commercial or financial  
372 relationships that could be construed as a potential conflict of interest.

#### ACKNOWLEDGEMENT

373 We thank Dr. Anthony Ritaccio for patient interactions, Dr. Aysegul Gunduz for help with data recording  
374 and Dr. Cuntai Guan for valuable discussions.

375 *Funding:* This work was supported by the NIH (EB00856, EB006356 and EB018783), the US  
376 Army Research Office (W911NF-08-1-0216, W911NF-12-1-0109, W911NF-14-1-0440) and Fondazione  
377 Neurone, and received support by the International Excellence Fund of Karlsruhe Institute of Technology.  
378 We acknowledge support by Deutsche Forschungsgemeinschaft and Open Access Publishing Fund of  
379 Karlsruhe Institute of Technology.

#### REFERENCES

- 380 Blakely, T., Miller, K. J., Rao, R. P., Holmes, M. D., and Ojemann, J. G. (2008), Localization  
381 and classification of phonemes using high spatial resolution electrocorticography (ECoG) grids,  
382 in *Engineering in Medicine and Biology Society*, 2008. EMBS 2008. 30th Annual International  
383 Conference of the IEEE (IEEE), 4964–4967
- 384 Bouchard, K. and Chang, E. (2014), Neural decoding of spoken vowels from human sensory-motor cortex  
385 with high-density electrocorticography, in *Engineering in Medicine and Biology Society*, 2014. EMBS  
386 2014. 36th Annual International Conference of the IEEE (IEEE)
- 387 Bouchard, K. E., Mesgarani, N., Johnson, K., and Chang, E. F. (2013), Functional organization of human  
388 sensorimotor cortex for speech articulation, *Nature*, 495, 7441, 327–332
- 389 Brumberg, J. S., Wright, E. J., Andreasen, D. S., Guenther, F. H., and Kennedy, P. R. (2011), Classification  
390 of intended phoneme production from chronic intracortical microelectrode recordings in speech-motor  
391 cortex, *Frontiers in neuroscience*, 5

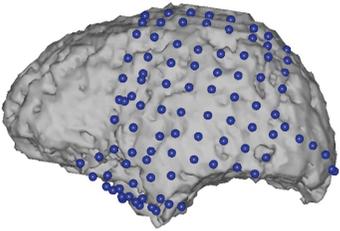
- 392 Canolty, R. T., Soltani, M., Dalal, S. S., Edwards, E., Dronkers, N. F., Nagarajan, S. S., et al. (2007),  
393 Spatiotemporal dynamics of word processing in the human brain, *Frontiers in neuroscience*, 1, 14
- 394 Chang, E. F., Rieger, J. W., Johnson, K., Berger, M. S., Barbaro, N. M., and Knight, R. T. (2010),  
395 Categorical speech representation in human superior temporal gyrus, *Nature neuroscience*, 13, 11,  
396 1428–1432
- 397 Crane, W., Gilbert, S., John, McConnell, W., Tenniel, S., John, Weir, H., and Zwecker, J. B. (1867),  
398 Mother Gooses Nursery Rhymes. A Collection of Alphabets, Rhymes, Tales and Jingles (London:  
399 George Routledge and Sons)
- 400 Crone, N., Hao, L., Hart, J., Boatman, D., Lesser, R., Irizarry, R., et al. (2001a), Electrographic  
401 gamma activity during word production in spoken and sign language, *Neurology*, 57, 11, 2045–2053
- 402 Crone, N. E., Boatman, D., Gordon, B., and Hao, L. (2001b), Induced electrographic gamma  
403 activity during auditory perception, *Clinical Neurophysiology*, 112, 4, 565–582
- 404 Deng, S., Srinivasan, R., Lappas, T., and D’Zmura, M. (2010), Eeg classification of imagined syllable  
405 rhythm using hilbert spectrum methods, *Journal of neural engineering*, 7, 4, 046006
- 406 Farwell, L. A. and Donchin, E. (1988), Talking off the top of your head: toward a mental prosthesis  
407 utilizing event-related brain potentials, *Electroencephalography and clinical Neurophysiology*, 70, 6,  
408 510–523
- 409 Formisano, E., De Martino, F., Bonte, M., and Goebel, R. (2008), ” who” is saying” what”? brain-based  
410 decoding of human voice and speech, *Science*, 322, 5903, 970–973
- 411 Fukuda, M., Rothmel, R., Juhász, C., Nishida, M., Sood, S., and Asano, E. (2010), Cortical gamma-  
412 oscillations modulated by listening and overt repetition of phonemes, *Neuroimage*, 49, 3, 2735–2745
- 413 Gales, M. and Young, S. (2008), The application of hidden markov models in speech recognition,  
414 *Foundations and Trends in Signal Processing*, 1, 3, 195–304
- 415 Gales, M. J. (1998), Maximum likelihood linear transformations for hmm-based speech recognition,  
416 *Computer speech & language*, 12, 2, 75–98
- 417 Gasser, T., Bächer, P., and Möcks, J. (1982), Transformations towards the normal distribution of broad  
418 band spectral parameters of the eeg, *Electroencephalography and clinical neurophysiology*, 53, 1, 119–  
419 124
- 420 Guenther, F. H., Brumberg, J. S., Wright, E. J., Nieto-Castanon, A., Tourville, J. A., Panko, M., et al.  
421 (2009), A wireless brain-machine interface for real-time speech synthesis, *PLoS one*, 4, 12, e8218
- 422 Haeb-Umbach, R. and Ney, H. (1992), Linear discriminant analysis for improved large vocabulary  
423 continuous speech recognition, in Acoustics, Speech, and Signal Processing, 1992. ICASSP-92.,  
424 1992 IEEE International Conference on, volume 1, volume 1, 13–16 vol.1, doi:10.1109/ICASSP.1992.  
425 225984
- 426 Huang, X., Acero, A., and Hon, H.-W. (2001), Spoken Language Processing: A Guide to Theory,  
427 Algorithm and System Development (Prentice Hall PTR)
- 428 Jelinek, F. (1997), Statistical methods for speech recognition (MIT press)
- 429 Kellis, S., Miller, K., Thomson, K., Brown, R., House, P., and Greger, B. (2010), Decoding spoken  
430 words using local field potentials recorded from the cortical surface, *Journal of neural engineering*, 7,  
431 5, 056007
- 432 Kennedy, J. F. (1989), Inaugural Addresses of the Presidents of the United States. (Washington, DC)  
433 (Available online at: [www.bartleby.com/124/](http://www.bartleby.com/124/))
- 434 Kubanek, J., Brunner, P., Gunduz, A., Poeppel, D., and Schalk, G. (2013), The tracking of speech envelope  
435 in the human cortex, *PLoS one*, 8, 1, e53398
- 436 Kubanek, J. and Schalk, G. (2014), Neuralact: A tool to visualize electrocortical (ecog) activity on a  
437 three-dimensional model of the cortex, *Neuroinformatics*, 1–8
- 438 Lee, K.-F. (1990), Context-dependent phonetic hidden markov models for speaker-independent  
439 continuous speech recognition, *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 38, 4,  
440 599–609
- 441 Leuthardt, E. C., Gaona, C., Sharma, M., Szrama, N., Roland, J., Freudenberg, Z., et al. (2011a), Using  
442 the electrocortical speech network to control a brain–computer interface in humans, *Journal of*  
443 *neural engineering*, 8, 3, 036004

- 444 Leuthardt, E. C., Pei, X.-M., Breshears, J., Gaona, C., Sharma, M., Freudenberg, Z., et al. (2011b),  
445 Temporal evolution of gamma activity in human cortex during an overt and covert word repetition  
446 task., *Frontiers in human neuroscience*, 6, 99–99
- 447 Lotte, F., Brumberg, J. S., Brunner, P., Gunduz, A., Ritaccio, A. L., Guan, C., et al. (2015),  
448 Electrographic representations of segmental features in continuous speech, *Frontiers in Human  
449 Neuroscience*, 9, 97, doi:10.3389/fnhum.2015.00097
- 450 Martin, S., Brunner, P., Holdgraf, C., Heinze, H.-J., Crone, N. E., Rieger, J., et al. (2014),  
451 Decoding spectrotemporal features of overt and covert speech from the human cortex, *Frontiers in  
452 Neuroengineering*, 7, 14, doi:10.3389/fneng.2014.00014
- 453 McFarland, D. J., Miner, L. A., Vaughan, T. M., and Wolpaw, J. R. (2000), Mu and beta rhythm  
454 topographies during motor imagery and actual movements, *Brain topography*, 12, 3, 177–186
- 455 Mesgarani, N., Cheung, C., Johnson, K., and Chang, E. F. (2014), Phonetic feature encoding in human  
456 superior temporal gyrus, *Science*, 1245994
- 457 Miller, K. J., Leuthardt, E. C., Schalk, G., Rao, R. P., Anderson, N. R., Moran, D. W., et al. (2007),  
458 Spectral changes in cortical surface potentials during motor movement, *The Journal of neuroscience*,  
459 27, 9, 2424–2432
- 460 Mugler, E., Goldrick, M., and Slutzky, M. (2014a), Cortical encoding of phonemic context during  
461 word production, in Engineering in Medicine and Biology Society, 2014. EMBS 2014. 36th Annual  
462 International Conference of the IEEE (IEEE)
- 463 Mugler, E. M., Patton, J. L., Flint, R. D., Wright, Z. A., Schuele, S. U., Rosenow, J., et al. (2014b), Direct  
464 classification of all american english phonemes using signals from functional speech motor cortex,  
465 *Journal of Neural Engineering*, 11, 3, 035015
- 466 Pasley, B. N., David, S. V., Mesgarani, N., Flinker, A., Shamma, S. A., Crone, N. E., et al. (2012),  
467 Reconstructing speech from human auditory cortex, *PLoS biology*, 10, 1, e1001251
- 468 Pei, X., Barbour, D. L., Leuthardt, E. C., and Schalk, G. (2011a), Decoding vowels and consonants  
469 in spoken and imagined words using electrocorticographic signals in humans, *Journal of neural  
470 engineering*, 8, 4, 046028
- 471 Pei, X., Leuthardt, E. C., Gaona, C. M., Brunner, P., Wolpaw, J. R., and Schalk, G. (2011b),  
472 Spatiotemporal dynamics of electrocorticographic high gamma activity during overt and covert word  
473 repetition, *Neuroimage*, 54, 4, 2960–2972
- 474 Potes, C., Gunduz, A., Brunner, P., and Schalk, G. (2012), Dynamics of electrocorticographic (ecog)  
475 activity in human temporal and frontal cortical areas during music listening, *NeuroImage*, 61, 4, 841 –  
476 848, doi:http://dx.doi.org/10.1016/j.neuroimage.2012.04.022
- 477 Pulvermüller, F., Huss, M., Kherif, F., del Prado Martin, F. M., Hauk, O., and Shtyrov, Y. (2006), Motor  
478 cortex maps articulatory features of speech sounds, *Proceedings of the National Academy of Sciences*,  
479 103, 20, 7865–7870
- 480 Rabiner, L. (1989), A tutorial on hidden markov models and selected applications in speech recognition,  
481 *Proceedings of the IEEE*, 77, 2, 257–286
- 482 Roy, E. and Basler, P. (1955), The gettysburg address, in *The Collected Works of Abraham Lincoln* (New  
483 Brunswick, NJ: Rutgers University Press)
- 484 Sahin, N. T., Pinker, S., Cash, S. S., Schomer, D., and Halgren, E. (2009), Sequential processing of  
485 lexical, grammatical, and phonological information within brocas area, *Science*, 326, 5951, 445–449
- 486 Schalk, G., McFarland, D. J., Hinterberger, T., Birbaumer, N., and Wolpaw, J. R. (2004), Bci2000: a  
487 general-purpose brain-computer interface (bci) system, *Biomedical Engineering, IEEE Transactions  
488 on*, 51, 6, 1034–1043
- 489 Schultz, T. and Kirchhoff, K. (2006), *Multilingual Speech Processing* (Elsevier, Academic Press)
- 490 Stolcke, A. (2002), Srilman extensible language modeling toolkit, in *Proceedings of the 7th International  
491 Conference on Spoken Language Processing (ICSLP 2002)*
- 492 Sutter, E. E. (1992), The brain response interface: communication through visually-induced electrical  
493 brain responses, *Journal of Microcomputer Applications*, 15, 1, 31–45
- 494 Talairach, J. and Tournoux, P. (1988), *Co-planar stereotaxic atlas of the human brain. 3-Dimensional  
495 proportional system: an approach to cerebral imaging* (Thieme)

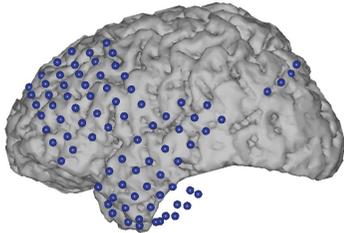
- 496 Telaar, D., Wand, M., Gehrig, D., Putze, F., Amma, C., Heger, D., et al. (2014), BioKIT - real-  
497 time decoder for biosignal processing, in The 15th Annual Conference of the International Speech  
498 Communication Association (Interspeech 2014)
- 499 Towle, V. L., Yoon, H.-A., Castelle, M., Edgar, J. C., Biassou, N. M., Frim, D. M., et al. (2008), Ecog  
500 gamma activity during a language task: differentiating expressive and receptive speech areas, *Brain*,  
501 131, 8, 2013–2027
- 502 unknown (2009), "Traitor among us" and "Split Feelings" (available on <https://www.fanfiction.net/>)
- 504 Wolpaw, J. R., Birbaumer, N., McFarland, D. J., Pfurtscheller, G., and Vaughan, T. M. (2002), Brain-  
505 computer interfaces for communication and control, *Clinical neurophysiology*, 113, 6, 767–791

Figure 1.JPEG

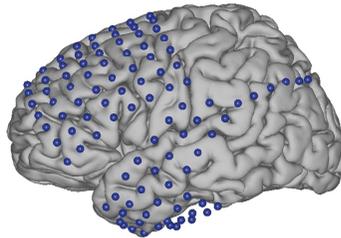
Subject 1  
29 y/o - F



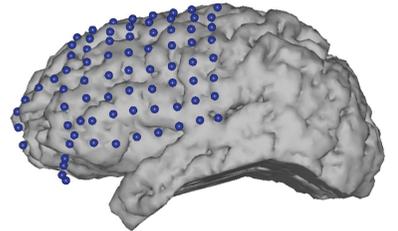
Subject 2  
30 y/o - M



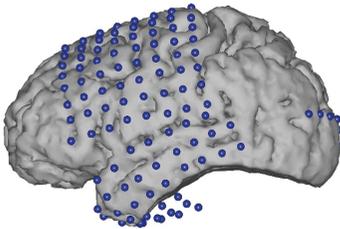
Subject 3  
29 y/o - F



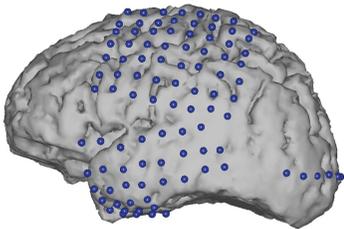
Subject 4  
18 y/o - M



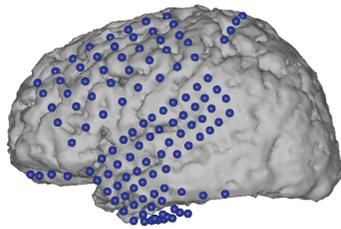
Subject 5  
26 y/o - F



Subject 6  
56 y/o - M



Subject 7  
29 y/o - F



Combined placement  
of subject 1, 2, 3, 5, 6, 7

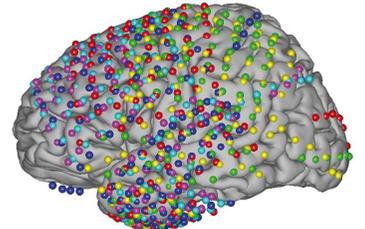


Figure 2.JPEG

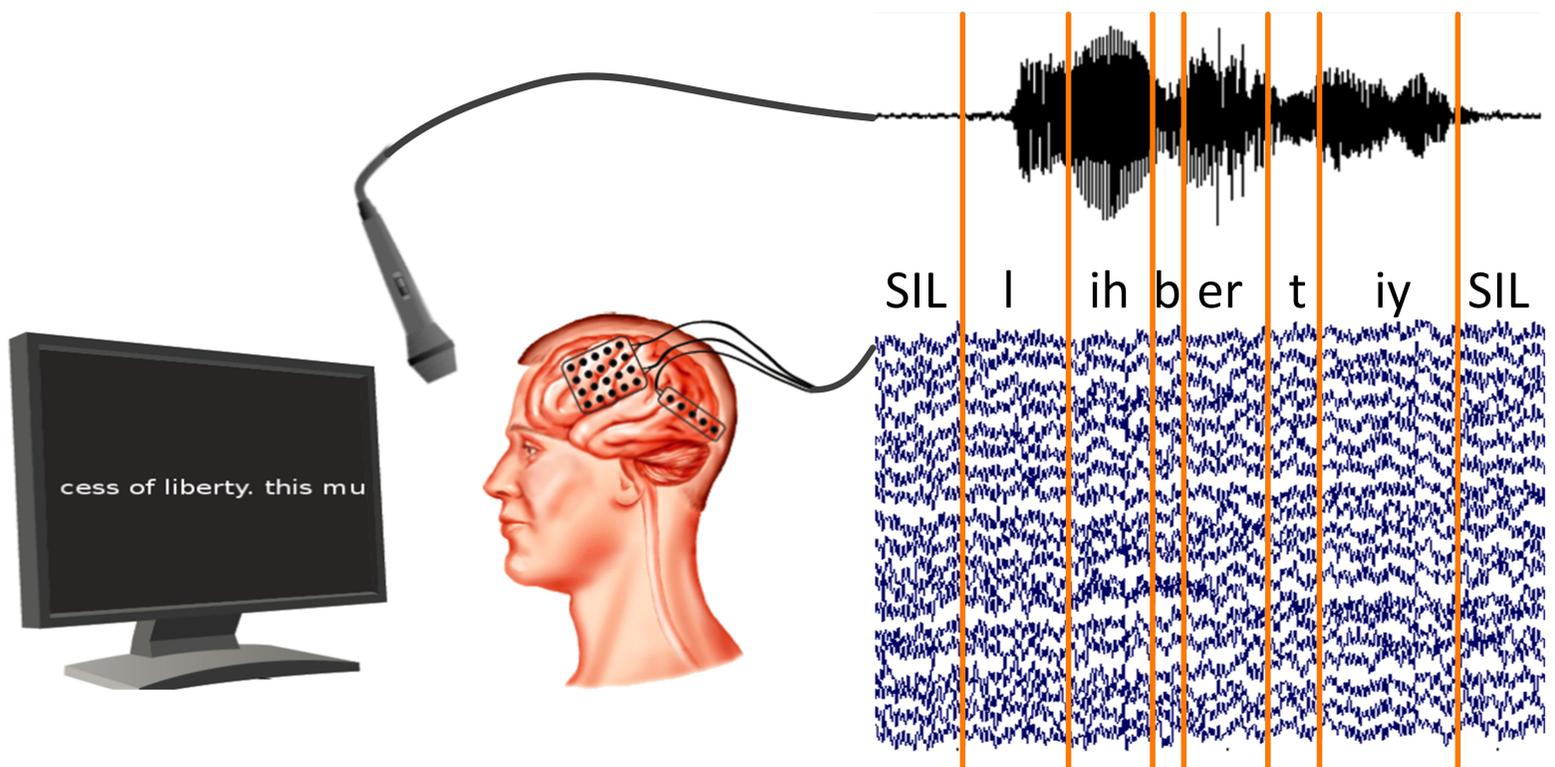


Figure 3.JPEG

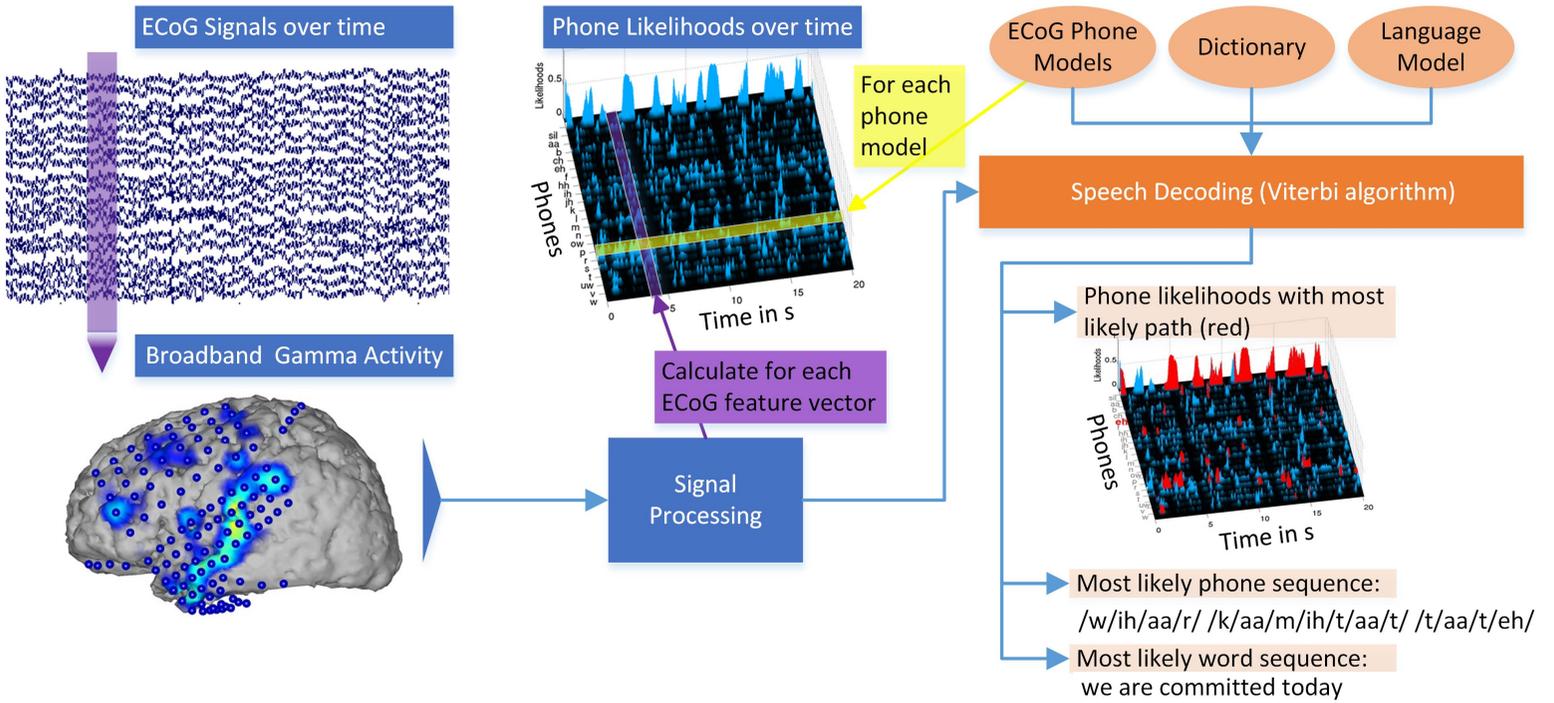


Figure 4.JPEG

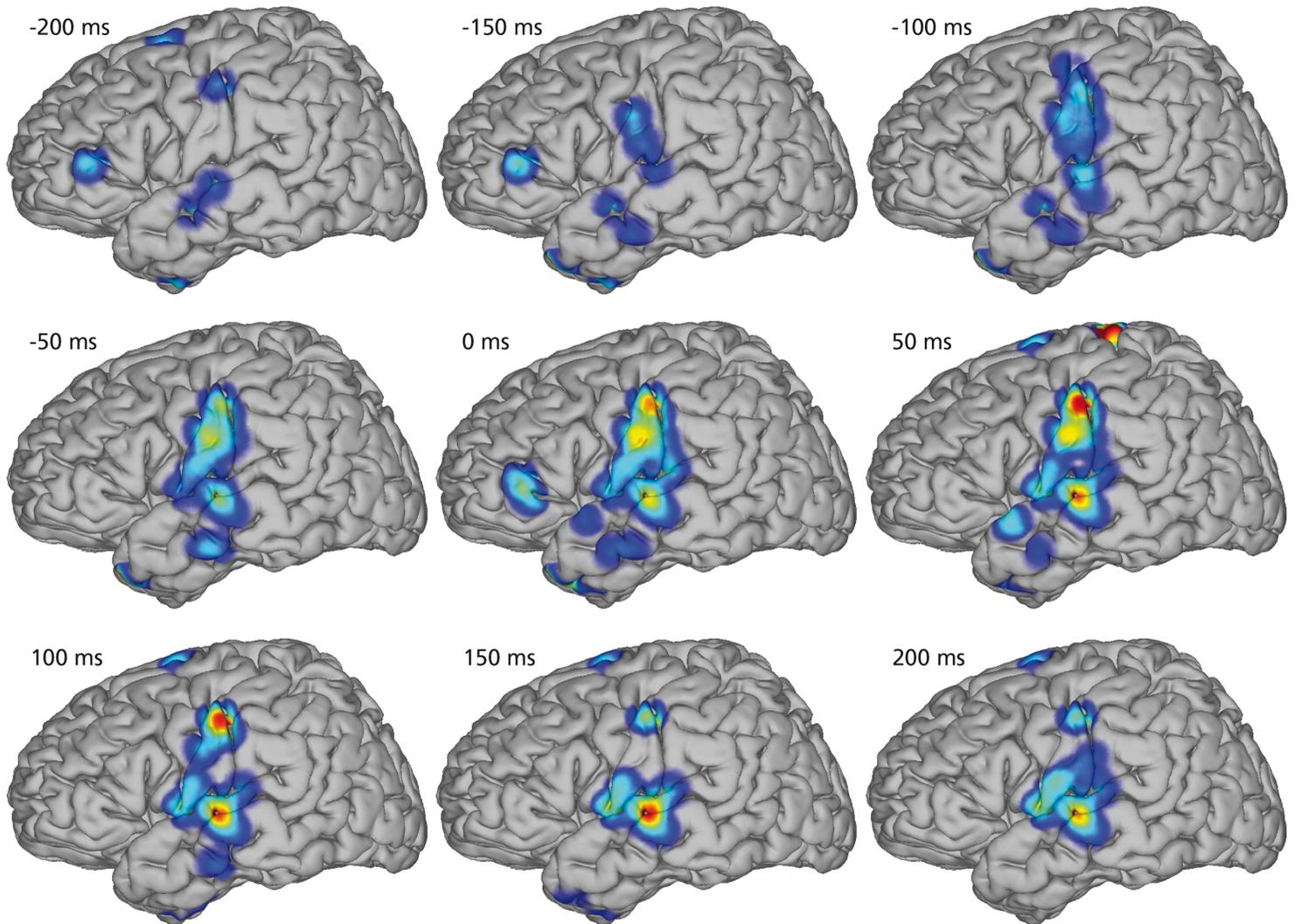


Figure 5.JPEG

