
Research Articles: Behavioral/Cognitive

Human sensorimotor cortex control of directly-measured vocal tract movements during vowel production

David F. Conant^{1,2,3}, Kristofer E. Bouchard^{4,5}, Matthew K. Leonard^{1,2} and Edward F. Chang^{1,2}

¹*Department of Neurological Surgery, University of California San Francisco (UCSF), San Francisco, California 94143*

²*Center for Integrative Neuroscience, UCSF, San Francisco, California 94143*

³*Graduate Program in Neuroscience, UCSF, San Francisco, California 94143*

⁴*Biological Systems and Engineering Division, Lawrence Berkeley National Laboratories (LBNL), Berkeley, CA 94720*

⁵*Helen Wills Neuroscience Institute, University of California, Berkeley (UCB), Berkeley, CA 94720*

DOI: 10.1523/JNEUROSCI.2382-17.2018

Received: 22 August 2017

Revised: 27 January 2018

Accepted: 29 January 2018

Published: 8 February 2018

Author contributions: D.F.C., K.E.B., and E.F.C. designed research; D.F.C. and K.E.B. collected the data; D.F.C. analyzed the data with assistance from K.E.B. and M.K.L.; D.F.C., K.E.B., M.K.L., and E.F.C wrote and edited the manuscript.; E.F.C. conceived and supervised the project.

Conflict of Interest: The authors declare no competing financial interests.

This work was supported by grants from the NIH (U01 NS098971 to E.F.C. and F32-DC013486 to M.K.L.). E.F.C. is a New York Stem Cell Foundation Robertson Investigator. This research was also supported by the New York Stem Cell Foundation, the Howard Hughes Medical Institute, the McKnight Foundation, The Shurl and Kay Curci Foundation, and The William K. Bowes Foundation. KEB was supported by Laboratory Directed Research and Development (LDRD) funding from Berkeley Lab, provided by the Director, Office of Science, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

Corresponding author: Edward Chang, University of California, San Francisco, Department of Neurological Surgery, 505 Parnassus Ave., M-779, San Francisco, CA 94143-0112, Email: edward.chang@ucsf.edu

Cite as: J. Neurosci ; 10.1523/JNEUROSCI.2382-17.2018

Alerts: Sign up at www.jneurosci.org/cgi/alerts to receive customized email alerts when the fully formatted version of this article is published.

Accepted manuscripts are peer-reviewed but have not been through the copyediting, formatting, or proofreading process.

Copyright © 2018 the authors

1 **Human sensorimotor cortex control of directly-measured vocal tract**
2 **movements during vowel production**

3
4 David F. Conant^{1,2,3}, Kristofer E. Bouchard^{4,5}, Matthew K. Leonard^{1,2}, Edward F.
5 Chang^{1,2*}
6

7
8 1. Department of Neurological Surgery, University of California San Francisco (UCSF), San Francisco,
9 California 94143

10 2. Center for Integrative Neuroscience, UCSF, San Francisco, California 94143

11 3. Graduate Program in Neuroscience, UCSF, San Francisco, California 94143

12 4. Biological Systems and Engineering Division, Lawrence Berkeley National Laboratories (LBNL),
13 Berkeley, CA 94720

14 5. Helen Wills Neuroscience Institute, University of California, Berkeley (UCB), Berkeley, CA 94720
15

16 *: Corresponding author

17 Edward Chang

18 University of California, San Francisco

19 Department of Neurological Surgery

20 505 Parnassus Ave., M-779

21 San Francisco, CA 94143-0112

22 Email: edward.chang@ucsf.edu
23
24

25 Author contributions: D.F.C., K.E.B., and E.F.C. designed research; D.F.C. and K.E.B.
26 collected the data; D.F.C. analyzed the data with assistance from K.E.B. and M.K.L.;
27 D.F.C., K.E.B., M.K.L., and E.F.C wrote and edited the manuscript.; E.F.C. conceived
28 and supervised the project.
29
30

31 This work was supported by grants from the NIH (U01 NS098971 to E.F.C. and F32-
32 DC013486 to M.K.L.). E.F.C. is a New York Stem Cell Foundation Robertson
33 Investigator. This research was also supported by the New York Stem Cell Foundation,
34 the Howard Hughes Medical Institute, the McKnight Foundation, The Shurl and Kay
35 Curci Foundation, and The William K. Bowes Foundation. KEB was supported by
36 Laboratory Directed Research and Development (LDRD) funding from Berkeley Lab,
37 provided by the Director, Office of Science, of the U.S. Department of Energy under
38 Contract No. DE-AC02-05CH11231.
39

40 The authors declare no competing financial interests.
41
42

43 8 Figures
44

45 Abstract: 227 Words

46 Introduction: 606 Words

47 Discussion: 1315 Words
48

49 Abstract

50 During speech production, we make vocal tract movements with remarkable
51 precision and speed. Our understanding of how the human brain achieves such proficient
52 control is limited, in part due to the challenge of simultaneously acquiring high-resolution
53 neural recordings and detailed vocal tract measurements. To overcome this challenge, we
54 combined ultrasound and video monitoring of the supralaryngeal articulators (lips, jaw
55 and tongue) with electrocorticographic (ECoG) recordings from the cortical surface of 4
56 subjects (three female, one male) to investigate how neural activity in the ventral sensory
57 motor cortex (vSMC) relates to measured articulator movement kinematics (position,
58 speed, velocity, acceleration) during the production of English vowels. We found that
59 high-gamma activity at many individual vSMC electrodes strongly encoded the
60 kinematics of one or more articulators, but less so for vowel formants and vowel identity.
61 Neural population decoding methods further revealed the structure of kinematic features
62 that distinguish vowels. Encoding and decoding of articulator kinematics was sparsely
63 distributed across time, and primarily occurred during the time of vowel onset and offset.
64 In contrast, encoding was low during the steady state portion of the vowel, despite
65 sustained neural activity at some electrodes. Significant representations were found for all
66 kinematic parameters, but speed was the most robust. These findings enabled by direct
67 vocal tract monitoring demonstrate novel insights to the representation of articulatory
68 kinematic parameters encoded in the vSMC during speech production.

69

70 Significance Statement

71 Speaking requires precise control and coordination of the vocal tract articulators
72 (lips, jaw, and tongue). Despite this impressive proficiency, our understanding of how the
73 brain achieves such control is rudimentary, in part because the movements themselves are
74 difficult to observe. By simultaneously measuring speech movements and the neural
75 activity that gives rise to them, we demonstrate how neural activity in sensorimotor
76 cortex produces complex, coordinated movements of the vocal tract.

77

78 Introduction

79 When we speak, we move the upper vocal tract articulators (lips, jaw, and tongue)
80 to produce vocal tract constrictions of air flow in precise, rapid, and complex ways.
81 These movements result in acoustic events that are highly distinguishable, maximizing
82 communicative utility. In spoken languages, vowels are a major category of speech
83 sounds. Despite their importance it is unknown how cortical activity patterns control the
84 vocal tract articulators to create vowels.

85 The ventral sensory-motor cortex (vSMC: pre-, post-, and sub-central gyri) is the
86 primary cortical area controlling the speech articulators (Petersen et al., 1988; Lotze et
87 al., 2000; Hesselmann et al., 2004; Brown et al., 2009; Kellis et al., 2010; Pei et al., 2011;
88 Bouchard et al., 2013; Bouchard and Chang, 2014; Mugler et al., 2014; Simonyan et al.,
89 2016; Kumar et al., 2016). Within vSMC, representations of vocal tract articulators are
90 coarsely somatotopically organized, with different neural populations in vSMC being
91 associated with specific articulators (Crone et al., 1998; Brown et al., 2009; Bouchard et
92 al., 2013; Mugler et al., 2014; Herff et al., 2015). However, our understanding of speech
93 motor control in vSMC is incomplete, due to challenges in simultaneously acquiring

94 neural and behavioral data with sufficient spatial and temporal resolution required to
95 determine the precise correspondence between vSMC activity and the movement of the
96 vocal tract articulators.

97 The precise movements (kinematics) of the articulators are challenging to
98 measure because many of the vocal tract movements are internal to the mouth and throat,
99 and therefore difficult to monitor externally, especially in the context of neural
100 recordings. As a result, previous studies have used the produced acoustics to infer which
101 articulators are involved, based on extensive linguistic descriptions of speech movements
102 for a given speech sound (Lotze et al., 2000; Crone et al., 2001; Brown et al., 2009;
103 Fukuda et al., 2010; Kellis et al., 2010; Pei et al., 2011; Leuthardt et al., 2011; Grabski et
104 al., 2012; Bouchard et al., 2013; Bouchard and Chang, 2014; Mugler et al., 2014; Herff et
105 al., 2015). Although it is possible to describe the movements of each articulator according
106 to phonetic labels derived from the acoustics, these behavioral descriptions cannot
107 provide exact characterizations of the changing positions of the articulators over time.
108 Moreover, there are many articulator configurations that can result in the same acoustics
109 (Atal et al., 1978; Maeda, 1990; Johnson et al., 1993; Gracco and Lofqvist, 1994;
110 Lofqvist and Gracco, 1999) and considerable across-speaker (Johnson et al., 1993) and
111 across-trial (Perkell and Nelson, 1985) variability in movements that give rise to a
112 particular speech sound. Thus, understanding how the brain produces complex sounds
113 like vowels requires determining how different kinematic parameters of articulatory
114 movements are controlled in vSMC during speech production.

115 To understand how vSMC neural activity controls precise articulator movements,
116 we have developed a system to simultaneously measure cortical activity using high-

117 resolution electrocorticography (ECoG) while directly monitoring the lips and jaw with a
118 camera, and the tongue with ultrasound. We previously detailed a technical description of
119 the methods (Bouchard et al., 2016). Here, we examined how vSMC generates articulator
120 kinematics, focusing on the production of American English vowels. We established that
121 articulator kinematics are more strongly represented in vSMC compared to acoustics. We
122 determined that specific kinematic parameters (position, speed, velocity, and
123 acceleration) are all represented, though articulator speed is represented most strongly.
124 Finally, we examined how distinct dynamics of neural activity are related to both
125 movement (from rest to target position) and maintenance of articulators (at target
126 position). By simultaneously measuring speech-related movements and the neural activity
127 generating them, we demonstrate how neural activity in sensorimotor cortex produce
128 complex, coordinated movements of the vocal tract.

129

130 Materials and Methods

131

132 *Experimental Design*

133

134 *Electrocorticography acquisition and signal processing*

135 Four human participants underwent chronic implantation of a high-density
136 subdural electrocorticographic array (ECoG) as part of the clinical treatment of epilepsy
137 (3 female right hemisphere, one male left hemisphere). All subjects were implanted with
138 256-channel grids over peri-Sylvian cortex (1.17mm diameter electrodes, 4mm pitch,
139 60x60mm coverage; Integra [Plainsboro NJ, USA]), referenced to scalp electrode. The

140 total number of vSMC electrodes for individual subjects ranged from 52 to 86 for a total
141 of 270. Cortical-surface electrical potentials were recorded with ECoG arrays and the
142 voltage time series from each electrode was inspected for artifacts or excessive noise.
143 Electrodes with excessive noise and time periods with artifacts were excluded from
144 analysis, and the raw ECoG activity was re-referenced to the common average. For each
145 channel, the time-varying analytic amplitude of the voltage signal in the high-gamma
146 (HG) range (70-150 Hz) was extracted using the Hilbert transform, according to
147 previously published procedures (Edwards et al., 2010). HG correlates well with multi-
148 unit firing (Ray and Maunsell, 2011), and has high spatial and temporal resolution
149 (Muller et al., 2016). The HG signal was down-sampled to 400 Hz for analysis and
150 plotting purposes. HG power was *z*-scored relative to activity recorded during periods of
151 silence during the same recording session. All analyses were limited to the ventral
152 sensory-motor cortex (vSMC), which was anatomically defined as the ventral portions of
153 the pre-central and post-central gyri, as well as the sub-central gyrus.

154 *Task*

155 Participants listened to audio recordings of nine English vowels (*a/æ/ʌ/ε/ɜ/i/i/ɔ/u*)
156 and were instructed to repeat each vowel. On each trial, to ensure they properly identified
157 the vowel, they first heard it in an /h-V-d/ context (e.g. 'hood'), and then they heard the
158 vowel in isolation. After a 1-1.5 sec delay, participants were presented with a visual cue
159 to produce the isolated vowel. They were not explicitly instructed to hold the vowel for a
160 specific amount of time. The median duration of production was 1.66 seconds (STD =
161 0.35 s). For each participant, between 15 and 30 repetitions of each vowel were collected
162 over the course of 3-6 recording sessions.

163

164 *Articulator Tracking*

165 We developed a system to record the movements of the main supra-laryngeal
166 articulators while participants performed the vowel production task (Figure 1A), the
167 details of which have been described previously (Bouchard et al., 2016). Briefly, to
168 capture the movement of the lips and jaw, a camera was placed in front of the
169 participant's mouth. The participant's lips were painted blue, and red dots were painted
170 on the tip of the nose and the chin to simplify the process of extracting the shape and
171 position of these articulators. The camera captured video at 30 frames per second. To
172 image the tongue, an ultrasound transducer was held firmly under the participant's chin
173 with the plane-of-view capturing the midline of the tongue. The ultrasound recorded
174 images at 30 frames per second, and the data were aligned to the lips/jaw video according
175 to the peak of the cross-correlation of the audio signals from each video. Using hue
176 thresholding, we extracted the lips and jaw automatically from these videos as binary
177 masks (Figure 1B). From these binary masks, we extracted the locations of the four
178 corners of the mouth (upper/lower lip, left/right corners) and the jaw. For the tongue, we
179 used EdgeTrak to extract 100 points of the mid-sagittal contour, which were then down
180 sampled to 3 points by taking the median x and y value for the front, middle, and back
181 thirds of the contour (Li et al., 2004). Since video and ultrasound were collected in
182 orthogonal spatial planes, x and y positions in the lips/jaw images reflect left/right and
183 top/bottom, whereas x and y positions in the tongue images reflect front/back and
184 top/bottom. To correct for differences in the relative position of the camera and
185 ultrasound transducer with respect to the participant, we referenced each articulatory

186 point to the neutral starting position at the beginning of each trial. From the measured
187 position of each articulatory feature (X) we also derived movement parameters including
188 velocity (X'), speed ($|X'|$), and acceleration (X'') of that articulator. We refer to these
189 parameters collectively as the articulator kinematics. While the lips and jaw were both
190 included in all analyses, we found that lip opening and jaw height were correlated for this
191 vowel production task (cross-subject average correlation: $r = 0.73 \pm 0.12$). Therefore, to
192 simplify visualizations we only show results for the lips.

193

194 *Statistical Analysis*

195

196 *Acoustic feature extraction*

197 Speech sounds were recorded using a Sennheiser microphone placed in front of
198 the participant's mouth and recorded at 22 kHz (Figure 1A). The speech signal was
199 transcribed offline using Praat (Boersma, 2001). For each vowel, we extracted the
200 formants (F_1 - F_4) using an inverse filter method (Watanabe, 2001; Ueda et al., 2007;
201 Bouchard et al., 2016).

202

203 *Trial duration standardization*

204 To standardize the durations of the vowels across trials and participants, we
205 linearly resampled each trial to be the median duration across vowels and subjects (1.66
206 seconds). Behavior and neural signals changed with rapid and stereotyped dynamics
207 around onset and offset; resampling the entire trial would systematically change those
208 dynamics based on vowel duration. Therefore, to preserve onset and offset dynamics, we

209 only resampled data in the time window from 250ms after the onset of the acoustics to
210 250ms before the offset: corresponding to the steady-state hold. Trials with durations less
211 than half or greater than twice the median were excluded from analysis (26 in total across
212 all subjects). Final analyses utilized an average of 15.3+/- 5.7 trials per vowel per subject.

213

214 *Permutation tests*

215 To evaluate statistical significance in each analysis, we used permutation tests. A
216 permutation distribution for a given model was constructed by randomly permuting the
217 trial labels of the observed data, and then training and testing the model using this
218 shuffled data. This process was repeated 500 times, and the performance of these shuffled
219 models comprised the permutation distribution. A model was considered significant if its
220 performance on test data was greater than the 99th percentile of its corresponding
221 permutation distribution. For the correlations in Figure 2D, we test if $|r| > 99^{\text{th}}$ percentile
222 of $|r_{\text{null}}|$.

223

224 *Correlations with articulatory position*

225 To evaluate the relationship between vSMC HG activity and individual
226 articulators, we correlated HG activity at individual electrodes with the measured trial-to-
227 trial position for each articulatory feature. HG activity averaged over a 200ms window
228 centered at acoustic onset was correlated with the mean position of each articulator taken
229 from a 200ms window centered at the midpoint of the vowel. Electrodes were labeled
230 according to whether they had significant correlations with zero, one, or multiple

231 articulator positions. Electrodes in 2D are examples of electrodes with significant
232 correlations with only one articulator.

233

234 *Encoding of kinematics, formants, and vowel categories*

235 We compared the representation of articulator kinematics, vowel formants, and
236 vowel category at each electrode using L_1 -regularized linear regression (Lasso). These
237 models predict HG activity at each time point 500ms before acoustic onset to 500ms after
238 acoustic offset from a sparse linear combination the behavior:

239 (1)
$$HG_e = \sum_{i=1}^n \beta_i$$

240 Where HG_e is the HG power at a given electrode, β are the linear weights that describe
241 the mapping, and i is a vowel category (n=9), vowel formant feature (n=10), articulator
242 kinematic feature (n=40), or all feature sets jointly (n=59). The formant features were F1-
243 F4, as well as all pairwise ratios of F1-F4. The articulator kinematic parameters were
244 position, speed, velocity, and acceleration for lip opening, lip width, jaw height, and the
245 front, middle, and back tongue. Vowel identity was parameterized as 9 binary vectors
246 corresponding to the vowel being produced during vocalization. Formant, articulatory,
247 and vowel identity features were lagged +100ms relative to HG, corresponding to the
248 causal direction of neural activity generating behavior. This lag was determined
249 empirically by optimizing model performance over a range of lag values (-500ms to
250 +500ms).

251 To train and test linear models, we used L_1 -regularized linear regression in a
252 leave-one-trial-out cross-validation procedure. We calculated the correlation between the
253 observed and predicted HG values, averaged across cross-validations. Electrodes were

254 included in visualizations and summary statistics only if their performance passed the
255 permutation test described above for at least one of these models (i.e., formants,
256 kinematics, vowel identity, or combined). To compare models with different numbers of
257 parameters, we calculated the adjusted R^2 :

$$258 \quad (2) \quad R_{adj}^2 = R^2 - (1 - R^2) \frac{p}{n-p-1}$$

259 Where R^2 is the unadjusted coefficient of determination of the model, n is the number of
260 observations the model was trained on, and p is the number of parameters.

261

262 *Organization of vowels in behavioral and neural spaces*

263 To examine the similarity of vowels in behavioral and neural representation
264 spaces, we used multi-dimensional scaling (MDS). MDS provides a low-dimensional
265 projection of the data that preserves the relative distances (similarities) between points in
266 a higher-dimensional space. For each feature set (formants, articulator position, and
267 neural) we extracted the median value for each vowel from a 200ms window centered at
268 the midpoint of the vowel (formants and articulator position) or the onset (neural), and
269 then z-scored that value across vowels. We applied MDS to the distance matrix computed
270 on these measurements for each feature set separately. To measure the differences in the
271 organization of vowels between the formant, articulator, and neural spaces, we calculated
272 the pairwise distances between the vowels in each space. We quantified the similarity
273 between the neural and kinematic or formant spaces by calculating a bootstrapped
274 correlation between the pairwise distances for each feature set. We performed
275 agglomerative hierarchical clustering on the pairwise distances to visually organize the
276 results.

277

278 *Encoding of kinematic parameters across time*

279 To assess the relative encoding of different kinematic parameters, we used the
280 measured position of each articulator on each trial (X) to derive the velocity (X'), speed
281 ($|X'|$), and acceleration (X'') of that articulator on that trial. To examine the encoding of
282 these parameters independent of one another, we removed the shared variance between
283 these parameters using semi-partial correlation. For each time point we first used linear
284 regression to predict the values of one kinematic parameter, y , from a linear combination
285 of the remaining 3 parameters, X :

$$286 \quad (3) \quad \hat{y} = \beta X$$

287 Where β are the weights that describe the linear relationship, and \hat{y} is the model's
288 prediction of that kinematic parameter. We then calculated the linearly independent
289 component of the kinematic parameter, y_{idp} , by subtracting predicted parameter values
290 from the observed:

$$291 \quad (4) \quad y_{idp} = y - \hat{y}$$

292 We then used L_1 -regularized linear encoding models to predict HG activity from the
293 kinematic parameters (position, speed, velocity, and acceleration) of the lips, jaw, and
294 tongue. However, instead of including the entire trial time-course in each model, we
295 trained and tested models within 100ms non-overlapping windows that tiled the trial.
296 Articulator kinematics were lagged +100ms relative to HG to evaluate the causal nature
297 of neural activity on behavior. Models were trained and tested independently for each
298 time window. Performance was measured by the correlation between the observed and
299 predicted HG values, averaged across cross-validations. Electrodes were included in

300 visualizations and summary statistics only if their performance passed the permutation
301 test described above for at least 3 contiguous time windows at any point in the trial.

302

303 *Decoding of kinematic parameters from vSMC HG*

304 To determine the degree to which variations in individual articulatory parameters
305 could be predicted from vSMC population activity, we used linear decoding models.
306 Similar to the encoding models above, we built L_1 -regularized linear models to relate
307 vSMC HG to articulator kinematics within 100ms time windows that tile the trial.
308 However, instead of predicting the HG activity at a single electrode from a combination
309 of all articulator parameters, we predicted the trial-to-trial variance of each articulator
310 parameter from a combination of all vSMC electrodes. As with the encoding models,
311 articulator kinematic features were lagged 100ms relative to vSMC HG, and models were
312 trained and tested independently for each time window using a leave-one-out cross-
313 validation procedure. The resulting models thus express how well the vSMC population
314 can predict each resulting articulator kinematic feature as a function of time within the
315 vocalization. To simplify visualization, we averaged performance across subjects and
316 similar articulators.

317

318 *Description of vSMC HG dynamics*

319 To characterize the major physiological response types in HG dynamics, we used
320 non-negative matrix factorization (NMF). NMF is a dimensionality reduction technique
321 that extracts a predetermined number (i.e., rank, k) of bases ($B \in \mathbb{R}^{m \times k}$) and weights
322 ($W \in \mathbb{R}^{n \times k}$) that linearly combine to reconstruct the non-negative data ($A \in \mathbb{R}^{m \times n}$), such

323 that $k < \min(n,m)$ under the constraint that both the bases and weights are strictly non-
 324 negative:

$$325 \quad (5) \quad A \approx BW^T; B, W \geq 0$$

326 The solutions B and W are found by solving the (bi-convex) constrained optimization
 327 problem:

$$328 \quad (6) \quad \hat{B}, \hat{W} = \min_{B,W} \frac{1}{2} \|A - BW^T\|_F^2; \text{ s.t. } B, W \geq 0$$

329 NMF is particularly useful for decomposing data into ‘parts’ that have interpretable
 330 meanings (e.g., transient vs. sustained response types) (Lee and Seung, 1999; Donoho
 331 and Stodden, 2004; Bouchard et al., 2016; Leonard et al., 2016). The HG activity for each
 332 vSMC electrode across all participants was averaged across trials, offset by the minimum
 333 value (such that all values were positive), and NMF was applied to the matrix of time
 334 courses x electrodes. To determine a parsimonious number of bases, we calculated the
 335 reconstruction error when projecting the data onto the NMF bases:

$$336 \quad (7) \quad err = \frac{1}{2} \|A - BW^T\|_F^2$$

337 We then found the number of bases (i.e., rank k) beyond which reconstruction error only
 338 marginally reduced (i.e., the elbow of the curve): five bases were used. The first two
 339 bases resembled the transient and sustained activity observed in Figure 7A. Electrodes
 340 with sustained activity were defined as those that had weighting for basis 1 greater than
 341 for basis 2. The width (HG_w) of the HG activity for sustained electrodes was derived as
 342 follows:

$$343 \quad (8) \quad HG_w = \operatorname{argmin} \int (HG_{e,t} - \overline{HG_{e,t}}) - \operatorname{argmax} \int (HG_{e,t} - \overline{HG_{e,t}})$$

344 Where $HG_{e,t}$ is the HG activity at given sustained electrode for a given trial. This
 345 measure was calculated for each sustained vSMC electrode, for each trial. We assessed

346 spatial organization by measuring the Euclidean distance between electrodes organized
347 according to their maximum NMF weight (i.e., transient or sustained). We compared
348 distributions of intra-parameter distances and cross-parameter distances to randomized
349 distributions derived by shuffling the labeling of the electrodes. If the HG dynamic
350 variability across vSMC is spatially organized, the distribution of intra-parameter and
351 cross-parameter distances should differ from the distributions of the random distributions.

352

353 Results

354 Participants produced nine English vowels in isolation (*a/æ/ʌ/ε/ɜ/ɪ/i/ʊ/u*) (e.g., the
355 vowels pronounced as in the following set of words: “calm”, “cat”, ”send”, ”fun”,
356 ”heard”, ”sit”, ”need”, ”should”, ”boot”) while neural activity was recorded from ventral
357 sensorimotor cortex (vSMC) and the movements of the supra-laryngeal articulators were
358 monitored. These vowels densely cover both the acoustic and kinematic space of all
359 American English vowels, and are a basic and essential component of all languages. We
360 specifically studied vowels in isolation for several reasons. First, the associated
361 movements of the speech articulators consist of a single displacement from rest, to the
362 target position, and back to rest. This simplicity provides the opportunity to study isolated
363 movements of the speech articulators free from the context of surrounding phonemes.
364 The task was also designed to minimize variability in the lower vocal tract (e.g. larynx),
365 which we did not explicitly monitor. Using the recorded acoustics, we verified that
366 subjects produced the vowels with little trial-to-trial variation in either pitch or intensity.
367 Additionally, the movements occur at distinct epochs, allowing us to resolve the neural
368 representation of the movement to the target, from the maintenance of that target, from
369 the return to the resting configuration.

370

371 *Articulator tracking during vowel production*

372 We simultaneously tracked the movements (Figure 1A) of the major
373 supralaryngeal articulators (i.e., lips, jaw, and tongue; Figure 1B) while recording neural
374 activity directly from the cortical surface (Figure 1C; see Methods, Bouchard et al.,
375 2016). We first verified that by extracting the positions of the articulators, we observed
376 characteristic vocal tract configurations that reflect distinct vowels. For example, the
377 vowel /a/ is characterized by lowering the front tongue, raising the back tongue, and
378 opening the lips, while the vowels /i/ and /u/ have different configurations (Figure 2A).
379 The measured articulatory movements captured these characteristics, and clearly
380 discriminated vowel categories (Figure 1B, 2B). We also used the produced acoustics as
381 a behavioral measure of vowel discriminability. By extracting the formants from the
382 acoustic signal, we observed distinct relative patterns of acoustic power for different
383 vowels. For example, /a/ is characterized by high F1 and low F2, whereas /i/ and /u/ have
384 different formant profiles (Figure 2C).

385

386 *Representation of articulator kinematics in ventral sensory-motor cortex*

387 These descriptions demonstrate that both articulator kinematics and acoustic
388 formants provide rich descriptions of the same behavior. However, although kinematics
389 and acoustics are causally related, their relationship is not one-to-one (Atal et al., 1978;
390 Maeda, 1990; Johnson et al., 1993; Gracco and Lofqvist, 1994; Lofqvist and Gracco,
391 1999), nor are they perfectly correlated (in the present data set, $\rho_{\text{kin,acoust}}=0.53 \pm 0.17$).
392 For example, producing the vowel /uw/ (“hoot”) involves raising the back of the tongue

393 towards the soft palate while rounding the lips. However, those movements can be
394 compensatory. The vowel /u/ can be produced with less pronounced lip movements
395 accompanied by greater tongue movements, or vice-versa (Perkell et al., 1993).
396 Therefore, we asked whether articulator kinematics or acoustic formants are the
397 behavioral characterization of vowels represented in vSMC.

398 First, we quantified how well the positions of speech articulators or vowel
399 formants explain the variance of HG at individual vSMC electrodes (i.e., encoding
400 strength). We recorded cortical electrical potentials from a total of 270 electrodes from
401 the surface of vSMC across 4 subjects (Figure 1C). The high-gamma (HG) activity at
402 many vSMC electrodes was elevated above baseline during the speech movements, and
403 was significantly correlated with the trial-to-trial position of the speech articulators
404 (Figure 2D). We observed a clear relationship between articulator position and HG
405 activity. For illustration, we identified representative electrodes where activity was most
406 correlated with a single articulator. For example, the HG activity of electrode 1 at the
407 time of vowel onset was significantly correlated with only the back tongue. Likewise,
408 electrode 2 showed greater activity for higher front tongue positions. Electrode 3 was
409 correlated with the opening of the lips. To examine whether HG activity at these
410 electrodes was similarly correlated with the produced acoustics, we binned the activity by
411 formant values (Figure 2E). We observed weaker correlations with formants compared to
412 articulator position, demonstrating more robust encoding of articulatory representations.

413 We were specifically interested in whether vSMC activity is best explained by
414 articulator kinematics, vowel formants, or vowel identity. We used linear encoding
415 models to predict neural activity from kinematic or acoustic features, or the vowel

416 identity (see Methods). Across electrodes, we found that articulator kinematics provided
417 significantly better model fits compared to vowel formants ($U = 5.3$, $p = 1.0e-8$;
418 Wilcoxon rank sum), or vowel identity ($U = 8.7$, $p = 4.1e-18$; Wilcoxon rank sum)
419 (Figure 3A). We used nested models to examine how much additional neural variance is
420 explained by predicting HG from both articulator kinematics and vowel formants. We
421 found that the joint model explained no more variance than the articulator kinematics
422 alone ($U = 0.4$, $p = 0.7$; Wilcoxon rank sum) suggesting that the performance of the
423 formant models was likely driven by variance shared with the kinematics. Therefore, we
424 find no evidence for encoding of vowel formants separate from their articulatory origin.
425 Furthermore, these results demonstrate that the production of distinct vowels is grounded
426 in direct control of articulator kinematics.

427 Across all vSMC electrodes, we found that 27% (74 out of 270) were significantly
428 correlated with movements of one or more articulators (correlations $> 99^{\text{th}}$ percentile of
429 permutation distributions). We observed a clear spatial organization to the articulator
430 correlations, with lips/jaw more dorsal than the tongue (Figure 3B), consistent with
431 previously-described somatotopy (Penfield and Roberts, 1959; Brown et al., 2009;
432 Bouchard et al., 2013; Huang et al., 2013; Conant et al., 2014). Within the ventral region,
433 we observed electrodes that more strongly reflected either the front or back of the tongue.
434 Both front and back tongue electrodes were distributed throughout the broader tongue
435 region. Finally, we observed 45 electrodes that had significant correlations with multiple
436 articulators, which were distributed throughout vSMC. Together, these results extend our
437 understanding of speech-motor cortex somatotopy by demonstrating that the dominant

438 encoding scheme in these neural populations reflects the specific movements of the
439 preferred articulators.

440

441 *Organization of vowels in vSMC population activity*

442 To understand how vSMC encoding of articulator kinematics contributes to our
443 ability to produce distinct vowels, we examined the organization of behavioral and neural
444 activity in relation to all nine vowels. In addition to articulator kinematic representations
445 at individual electrodes, population activity in vSMC may reflect the coordinated
446 movements of the vocal tract that produce vowels. Furthermore, because vowel formants
447 arise from the relative positions of multiple articulators, it may be the case that while
448 articulators are most strongly represented at single electrodes, population activity may
449 reflect a different, emergent representation. We examined the organization of speech
450 representations at the population level by comparing the relative distances of vowels in
451 acoustic, articulatory, and neural space. We performed multi-dimensional scaling (MDS)
452 on the vowel centroids (see Methods) measured by vowel formants, articulator position,
453 or vSMC neural activity across all participants. In this analysis, vowel tokens that are
454 near each other in MDS space have similar formant, kinematic, or HG values. Consistent
455 with previous behavioral and linguistic descriptions of vowel production, the formant and
456 kinematic MDS projections replicate the classic vowel space ‘trapezoid’ (Figure 4A-C)
457 (Hillenbrand et al., 1995; Ladefoged and Johnson, 2011). The HG neural MDS projection
458 also closely resembled the acoustic/kinematic organization of the vowels. For example,
459 the vowel /a/ (as in hall) is near the vowel /ʌ/ (as in hut), but far from /i/ (as in heat).

460 To characterize the difference in organization of vowels across these feature
461 spaces, we calculated the pairwise distances between the vowels in MDS space,
462 visualized as confusion matrices (Figure 4D-F, right). We additionally performed
463 hierarchical clustering of the pairwise distances and organized the confusion matrices by
464 the derived hierarchical organization. The pairwise distances and hierarchical clustering
465 reaffirm the classic vowel organization, but also highlight the specific differences
466 between the feature spaces. For example, /i/ is distant from the other vowels in the
467 formant space, but closer in the articulator and neural spaces. We found that the
468 organization of vowels in vSMC HG activity is significantly more correlated with the
469 organization of vowels in the articulator space compared to the acoustic space ($U = 9,5$,
470 $p=1.3e-21$, Wilcoxon rank sum), although both representations were significantly
471 correlated with vowel organization in the HG neural space (acoustic: $r=0.56$, $p=2.8e-4$,
472 kinematic: $r=0.73$, $p=5.9e-5$).

473

474 *Encoding of articulator kinematic parameters*

475 In the above analyses, we considered the joint encoding of multiple kinematic
476 parameters for individual articulators. However, it is unknown whether kinematic
477 encoding reflects particular aspects of the articulator movements. The movements of the
478 articulators can be described according to a variety of different kinematic parameters (e.g.
479 position, speed or velocity, acceleration, etc.). For each kinematic parameter, we used
480 L1-regularized encoding models to explain vSMC HG from the moment-to-moment
481 measurements of position, speed, velocity, and acceleration. Since all four kinematic
482 parameters are correlated with one another, we removed shared variance between the

483 parameters using semi-partial correlations in order to better interpret their relative
484 encoding performances.

485 We found electrodes that significantly encoded the trial-to-trial variability in
486 position (Figure 5A), speed (Figure 5B), velocity (Figure 5C), and acceleration (Figure
487 5D). Speed was the most robustly encoded parameter at most vSMC electrodes, with
488 significant encoding at more electrodes and a higher average correlation compared to the
489 other parameters ($U = 1720$ to 2735 , $p = 3.3e-9$ to $1.1e-14$; Wilcoxon rank sum; Figure
490 5E).

491 To understand the timing of kinematic parameter encoding throughout the
492 production of vowels, we also examined models that predicted HG neural activity from
493 the joint combination of all four parameters simultaneously. These models were
494 evaluated over a sliding 100ms window to characterize the kinematic parameter encoding
495 during different phases of the trial (i.e., movement initiation, target position, steady state
496 maintenance, and movement back to the starting position). We observed a peak in
497 encoding for most electrodes around the onset of the movement (91% of electrodes), with
498 some electrodes also showing a peak around the offset (9%; Figure. 5F). There was no
499 spatial organization associated with electrodes that specifically encoded particular
500 parameters (intra-parameter, $p = 0.31$; cross-parameter, $p = 0.08$; see Methods), nor was
501 there a significant relationship between electrodes that encoded specific kinematic
502 parameters and specific articulators ($\chi^2(9, N=155) = 9.26$, $p = 0.4$; Chi-square).
503 Strikingly, encoding during the steady state was near zero for all kinematic parameters.

504 To understand how these individual electrode kinematic representations relate to
505 the population representations of articulator kinematics and dynamics, we additionally

506 used L1-regularized decoding models to predict the articulator kinematics from the
507 population of vSMC HG electrodes. As with the encoding analyses, these models were
508 constructed from a small (100ms) sliding window of time, resulting in a description of
509 how much of the trial-to-trial variability of the articulator position (Figure 6A), speed
510 (Figure 6B), velocity (Figure 6C), and acceleration (Figure 6D) can be explained by
511 vSMC HG activity. The time course of decoding strength was similar to the encoding
512 models, with peaks around the onset and offset and near-zero values while the vowel was
513 being held. Across kinematic parameters, articulator speed was the best-predicted
514 parameter ($U = 2.6$ to 3.3 , $p = 8.1e-3$ to $9.0e-4$; Wilcoxon rank sum).

515 Together, these results demonstrate a strikingly sparse representation of kinematic
516 parameters across time, despite the fact that there continues to be trial-to-trial variability
517 in both kinematic and neural features throughout vowel (Figure 5). Only 56% of time
518 points had significant encoding performance at any electrode, and individual significant
519 electrodes had an average of 15% (± 1) significant time points. In particular, we did not
520 observe any electrodes that exhibited significant kinematic parameter encoding during the
521 steady state of the vowel.

522

523 *Onset vs. steady state HG activity and kinematic encoding*

524 The temporal sparsity of neural representations described above is particularly
525 notable given that many electrodes showed sustained HG activity during the steady state
526 portion of the vowel, independent of the particular articulatory movements that occurred
527 (Figure 7A). These electrodes contrast with other HG activity that is only transiently
528 increased around the onset and/or offset of the vowel (Figure 7A). To characterize these

529 response types, we used non-negative matrix factorization (NMF) to derive basis
530 functions that best describe vSMC HG temporal profiles across all electrodes (Hamilton
531 et al., 2017). Our motivation for using NMF was not to provide a complete description of
532 HG dynamics, but rather to provide an unsupervised method of quantifying the transient
533 and onset/offset responses across electrodes. We found that the first two bases (i.e., the
534 most important bases), captured the sustained and onset/offset response types we
535 observed qualitatively (Figure 7B). Organizing all vSMC electrodes by the degree to
536 which their activity is reconstructed by the first or second NMF bases (i.e., the NMF
537 weights), we observed a continuum of HG dynamics: some electrodes showed sustained
538 activity throughout the utterance, while others showed transient increases in activity only
539 at onset and offset of the utterance (Figure 7C). Some electrodes showed a combination
540 of sustained and transient components. There was no apparent spatial organization (intra-
541 response type, $p = 0.13$; cross-response type $p = 0.09$; see Methods) or relationship
542 between response type and the articulators represented at each electrode ($\chi^2(3, N=155) =$
543 2.3 , $p = 0.5$; Chi-square).

544 We separately considered electrodes that showed stronger weights for the
545 sustained NMF basis (basis 1 in Figure 7B). The average HG activity at these electrodes
546 indeed showed sustained activity throughout the vowel, however there was not a
547 concomitant sustained encoding of kinematic parameters (Figure 7D). This dissociation
548 between activity and encoding was apparent even at the single trial level (Figure 8A).
549 Thus, although some electrodes exhibit sustained HG activity throughout the production
550 of the vowel, there is not a systematic relationship between the trial-to-trial variability of
551 that activity and the kinematics of the articulators. Instead, encoding of kinematics at

552 electrodes with sustained activity was prevalent only around the onset and offset of
553 movement. We hypothesized that although activity during the steady state does not relate
554 to kinematic variability, it still reflects an important aspect of the task, namely the
555 duration of each utterance. Across sustained electrodes we found that the duration of the
556 HG timecourse was significantly correlated with the duration of the vocalization
557 (Spearman's $\rho = 0.61$, $p = 2e-153$; Figure 8B). Thus, at a minimum the sustained
558 activity was associated with vowel production.

559

560 Discussion

561 We report a detailed description of how activity in speech-motor cortex controls
562 the precise movements of the vocal tract articulators to produce vowels. By
563 simultaneously measuring the movements of the articulators, recording the acoustic
564 consequences of those movements, and recording the neural activity in vSMC, we are
565 able to establish that the dominant representation in vSMC is articulator kinematics. The
566 precise control of these movements allows speakers to create specific configurations of
567 the mouth, which lead to distinct categories of sounds.

568 Without simultaneous measurements of the articulators, previous studies of the
569 neural basis of speech production have relied on approximate, categorical phonemic-
570 based descriptions of speech behavior (Crone et al., 2001; Fukuda et al., 2010; Kellis et
571 al., 2010; Leuthardt et al., 2011; Pei et al., 2011; Bouchard et al., 2013; Bouchard and
572 Chang, 2014; Mugler et al., 2014; Herff et al., 2015). Although the produced acoustics
573 and categorical vowel descriptions reflect the ultimate (perceptual) outcome of vocal tract
574 movements, the many-to-one relationship between kinematics and vowels (Atal et al.,

575 1978; Maeda, 1990; Johnson et al., 1993; Gracco and Lofqvist, 1994; Lofqvist and
576 Gracco, 1999) means that it was not possible to understand the precise nature of the
577 neural representation in vSMC. Previous studies have implicated vSMC in articulator
578 kinematic control in several ways. First, stimulation to sites in vSMC elicits involuntary
579 activations of the orofacial muscles (Penfield and Boldrey, 1937; Huang et al., 2013).
580 Second, neurons in these and other sensorimotor regions are often tuned to movement
581 kinematics (Georgopoulos et al., 1986; Paninski et al., 2004; Arce et al., 2013). Finally,
582 the spatio-temporal patterns of HG activity in vSMC are consistent with the engagement
583 of the articulators (Bouchard et al., 2013). The present results confirm these
584 interpretations by showing directly that kinematic descriptions of speech behavior are
585 more closely related to neural activity compared to acoustic or categorical vowel
586 descriptions. Further, we find no evidence that vSMC activity encodes either produced
587 acoustics or vowel category distinct from their correlations with the articulator
588 kinematics. Crucially, we observed that this encoding scheme exists both at single
589 electrodes, and across the spatially distributed set of electrodes. For spatially distributed
590 activity patterns, we demonstrate the neural existence of the classic vowel ‘trapezoid’,
591 which has dominated linguistic descriptions of speech (Harshman et al., 1977; Alfonso
592 and Baer, 1982; Hillenbrand et al., 1995).

593 Furthermore, by characterizing the movements of the articulators according to a
594 variety of kinematic parameters (position, speed, velocity, and acceleration), we
595 demonstrated that neural activity encodes each of the examined parameters independent
596 of one another. Previous studies examining arm movements using analogous parameters
597 have also found significant encoding of these parameters (Georgopoulos et al., 1982,

598 1984; Ashe et al., 1994; Moran and Schwartz, 1999; Paninski et al., 2004). While we find
599 electrodes that significantly encode each parameter examined, speed is by far the most
600 robustly encoded parameter. Furthermore, the dominant kinematic parameter at
601 individual electrodes was not significantly related to the articulator representation of
602 those electrodes. The predominance of speed over other parameters is somewhat
603 surprising; previous studies of the single-unit representation of kinematic parameters
604 during arm reaching typically find that velocity and direction are the most commonly
605 encoded parameter (Moran and Schwartz, 1999). Similar results were also observed in a
606 recent ECoG study, which found that movement speed was predominately represented
607 during arm reaching in humans (Hammer et al., 2016). The predominance of speed
608 encoding was interpreted in the context of a model in which the summed activity of many
609 velocity-tuned neurons with random directional tuning resembles speed tuning. Thus it
610 may be the case that individual vSMC neurons are actually representing mostly velocity,
611 but the summed activity observed with ECoG electrodes reflects the magnitude of
612 movement without direction (i.e. speed).

613 By studying vowels, we were able to examine the dynamics of kinematic
614 encoding that are associated with movements to specific vocal tract articulators. We
615 found that articulator kinematics were encoded around the time of movement onset
616 and/or offset, but not while the vocal tract configuration was being held to maintain the
617 vowel. Encoding of articulator kinematics only during movement onset and offset
618 suggests that control of speech articulators is accomplished primarily through control of
619 changes to the plant, rather than moment-to-moment maintenance of the vocal tract
620 configuration. This is consistent with models of speech production that utilize changes to

621 the plant as the primary mechanism by which sensorimotor cortex receives input from,
622 and sends commands to, the vocal tract (Houde and Nagarajan, 2011; Tourville and
623 Guenther, 2011). Furthermore, these dynamics have been observed in studies with
624 analogous behavior from different body parts, including arm reaching. These studies have
625 found individual neurons in motor cortex that exhibit transient firing patterns, where
626 firing rates are high around movement onset and offset (Crammond and Kalaska, 1996;
627 Shalit et al., 2012; Arce et al., 2013; Shadmehr, 2017).

628 We also found a subset of electrodes that exhibited sustained neural activity
629 during the steady-state portion of the vowel which was not correlated with any measured
630 kinematic features. Instead, we found that the duration of the sustained activity correlated
631 well with trial-by-trial vowel length. At a minimum, this suggests that this sustained
632 activity co-varies with whether the subject is vocalizing. One possibility is that sustained
633 activity represents an articulatory parameter that has little variability in our task, such as
634 respiration. However, a more intriguing possibility is that sustained activity may
635 represent a non-specific signal for holding the vocal tract configuration, which does not
636 directly encode the articulatory kinematics like position. Such a signal combined with the
637 onset/offset encoding of kinematics may provide sufficient information for encoding the
638 observed behavior. Further studies utilizing tasks with more variability in manner of
639 articulation are necessary to resolve these possibilities.

640 It is important to emphasize that these analyses focus on the neural representation
641 of the supra-laryngeal articulators. While the movements of these articulators are critical
642 to the production of vowels, the lower vocal tract (e.g. larynx, pharynx, and diaphragm)
643 is also necessary to produce voiced sounds. It is likely that sub-regions of vSMC are

644 involved in the control of the lower vocal tract (Brown et al., 2008; Bouchard et al., 2013;
645 Conant et al., 2014), but the limitations of our vocal tract kinematic monitoring system
646 and the lack of across-trial variability in pitch and intensity preclude a detailed
647 examination of the representations of these articulators in the present experiment.

648 Further, we are not able to make evaluate whether the activity we observe is due
649 to feed-forward signals originating in vSMC, or sensory feedback signals. Our models
650 performed optimally at a neural-leading lag of approximately 100ms, implying that the
651 representations we observed were driven more by feed-forward activity. However, the
652 relatively simple movements examined here exhibit temporal auto-correlation, which
653 makes it difficult to dissociate feed-forward activity from feedback. Examining speech
654 tasks with faster, less stereotyped movements (e.g. naturally produced words or
655 sentences) would make it possible to disentangle feed-forward and feedback signals, and
656 is an interesting and important future direction (Chang et al., 2013; Greenlee et al., 2013;
657 Kingyon et al., 2015; Behroozmand et al., 2016; Li et al., 2016; Cao et al., 2017).

658 Finally, while we observed qualitatively similar results across patients regardless
659 of the hemisphere in which the electrodes were implanted, three of the four participants
660 had grid placements on the right (non-dominant) hemisphere. It is presently unknown
661 whether there are differences in the representation of articulator movements between left
662 and right hemisphere, and the results presented here may not fully address the extent to
663 which such differences exist.

664 We found that the representation of spoken vowels in vSMC is directly explained
665 by the movements of speech articulators. The encoding of multiple kinematic parameters
666 is present for the articulators, most prominently speed. Articulator kinematic encoding

667 was primarily observed at the onset and offset of vowel production and not while the
 668 vowel is being held. Together, these findings provide insight into how neural activity in
 669 primary sensorimotor cortex results in the precise production of vowels. Future work will
 670 address how these encoding properties operate in the context of natural continuous
 671 speech.

672

673 **References**

674

- 675 Alfonso PJ, Baer T (1982) Dynamics of Vowel Articulation. *Lang Speech* 25:151–173.
 676 Arce FI, Lee J-C, Ross CF, Sessle BJ, Hatsopoulos NG (2013) Directional information
 677 from neuronal ensembles in the primate orofacial sensorimotor cortex. *J*
 678 *Neurophysiol* 110:1357–1369.
 679 Ashe J, Georgopoulos AP, Medical VA (1994) Movement Parameters and Neural
 680 Activity in Motor Cortex and Area 5. *Cereb Cortex* 6:590–600.
 681 Atal BS, Chang JJ, Mathews M V, Tukey JW (1978) Inversion of articulatory-to-acoustic
 682 transformation in the vocal tract by a computer-sorting technique. *J Acoust Soc Am*
 683 63:1535–1553.
 684 Behroozmand R, Oya XH, Nourski K V, Kawasaki H, Larson CR, Brugge JF, Howard
 685 MA, Greenlee JDW (2016) Neural Correlates of Vocal Production and Motor
 686 Control in Human Heschl’s Gyrus. *J Neurosci* 36:2302–2315.
 687 Boersma P (2001) Praat, a system for doing phonetics by computer. *Glott Int* 5:341–345.
 688 Bouchard KE, Chang EF (2014) Control of Spoken Vowel Acoustics and the Influence of
 689 Phonetic Context in Human Speech Sensorimotor Cortex. *J Neurosci* 34:12662–
 690 12677.
 691 Bouchard KE, Conant DF, Anumanchipalli GK, Dichter B, Chaisanguanthum KS,
 692 Johnson K, Chang EF (2016) High-resolution, non-invasive imaging of upper vocal
 693 tract articulators compatible with human brain recordings. *PLoS One* 11:1–30.
 694 Bouchard KE, Mesgarani N, Johnson K, Chang EF (2013) Functional organization of
 695 human sensorimotor cortex for speech articulation. *Nature* 495:327–332.
 696 Brown S, Laird AR, Pfordresher PQ, Thelen SM (2009) The somatotopy of speech:
 697 Phonation and articulation in the human motor cortex. *Brain Cogn* 70:31–41.
 698 Cao L, Thut G, Gross J (2017) The role of brain oscillations in predicting self-generated
 699 sounds. *Neuroimage* 147:895–903.
 700 Chang EF, Niziolek C a, Knight RT, Nagarajan SS, Houde JF (2013) Human cortical
 701 sensorimotor network underlying feedback control of vocal pitch. *Proc Natl Acad*
 702 *Sci U S A* 110:2653–2658.
 703 Conant D, Bouchard KE, Chang EF (2014) Speech map in the human ventral sensory-
 704 motor cortex. *Curr Opin Neurobiol* 24:63–67.
 705 Crammond D, Kalaska JF (1996) Differential relation of discharge in primary motor
 706 cortex and premotor cortex to movements versus actively maintained postures

- 707 during a reaching task. *Exp brain Res* 108:45–61.
- 708 Crone NE, Hao L, Hart J, Boatman D, Lesser RP, Irizarry R, Gordon B (2001)
- 709 Electrographic gamma activity during word production in spoken and sign
- 710 language. *Neurology* 57:2045–2053.
- 711 Crone NE, Miglioretti DL, Gordon B, Lesser RP (1998) Functional mapping of human
- 712 sensorimotor cortex with electrocorticographic spectral analysis. II. Event-related
- 713 synchronization in the gamma band. *Brain* 121:2301–2315.
- 714 Donoho D, Stodden V (2004) When does non-negative matrix factorization give a correct
- 715 decomposition into parts? *Proc Adv Neural Inf Process Syst* 16:1141–1148.
- 716 Edwards E, Nagarajan SS, Dalal SS, Canolty RT, Kirsch HE, Barbaro NM, Knight RT
- 717 (2010) Spatiotemporal imaging of cortical activation during verb generation and
- 718 picture naming. *Neuroimage* 50:291–301.
- 719 Fukuda M, Rothermel R, Juhász C, Nishida M, Sood S, Asano E (2010) Cortical gamma-
- 720 oscillations modulated by listening and overt repetition of phonemes. *Neuroimage*
- 721 49:2735–2745.
- 722 Georgopoulos a P, Kalaska JF, Caminiti R, Massey JT (1982) On the relations between
- 723 the direction of two-dimensional arm movements and cell discharge in primate
- 724 motor cortex. *J Neurosci* 2:1527–1537.
- 725 Georgopoulos AP, Caminiti R, Kalaska JF (1984) Static Spatial Effects in Motor Cortex
- 726 and Area 5: Quantitative Relations in a Two-Dimensional Space. *Exp brain*
- 727 *Res*:446–454.
- 728 Georgopoulos AP, Schwarz AB, Ketner RE (1986) Neuronal population coding of
- 729 movement direction. *Science* (80-) 233:1416–1419.
- 730 Grabski K, Lamalle L, Vilain C, Schwartz J-L, Vallée N, Tropres I, Baciú M, Le Bas J-F,
- 731 Sato M (2012) Functional MRI assessment of orofacial articulators: neural correlates
- 732 of lip, jaw, larynx, and tongue movements. *Hum Brain Mapp* 33:2306–2321.
- 733 Gracco VL, Lofqvist A (1994) Speech motor coordination and control: evidence from lip,
- 734 jaw, and laryngeal movements. *J Neurosci* 14:6585–6597.
- 735 Greenlee JDW, Behroozmand R, Larson CR, Jackson AW, Chen F, Hansen DR, Oya H,
- 736 Kawasaki H, Howard M a (2013) Sensory-motor interactions for vocal pitch
- 737 monitoring in non-primary human auditory cortex. *PLoS One* 8:e60783.
- 738 Hamilton LS, Edwards E, Chang EF (2017) Parallel streams define the temporal
- 739 dynamics of speech processing across human auditory cortex. *bioRxiv*.
- 740 Hammer J, Pistohl T, Fischer J, Kr P, Tomá M, Marusic P, Schulze-bonhage A, Aertsen
- 741 A, Ball T (2016) Predominance of Movement Speed Over Direction in Neuronal
- 742 Population Signals of Motor Cortex : Intracranial EEG Data and A Simple
- 743 Explanatory Model. *Cereb Cortex* 26:2863–2881.
- 744 Harshman R, Ladefoged P, Goldstein L (1977) Factor analysis of tongue shapes. *J Acoust*
- 745 *Soc Am* 62:693–713.
- 746 Herff C, Heger D, De pesters A, Telaar D, Brunner P, Schalk G, Schultz T (2015) Brain-
- 747 to-text: Decoding spoken phrases from phone representations in the brain. *Front*
- 748 *Neuroeng* 8.
- 749 Hesslmann V, Sorger B, Lasek K, Guntinas-Lichius O, Krug B, Sturm V, Goebel R,
- 750 Lackner K (2004) Discriminating the cortical representation sites of tongue and up
- 751 movement by functional MRI. *Brain Topogr* 16:159–167.
- 752 Hillenbrand J, Getty L a., Clark MJ, Wheeler K (1995) Acoustic characteristics of

- 753 American English vowels. *J Acoust Soc Am* 97:3099–3111.
- 754 Houde JF, Nagarajan SS (2011) Speech production as state feedback control. *Front Hum*
755 *Neurosci* 5:1–14.
- 756 Huang CS, Sirisko MA, Hiraba H, Murray GM, Sessle J (2013) Organization of the
757 primate face motor cortex as revealed by intracortical microstimulation and
758 electrophysiological identification of afferent inputs and corticobulbar projections
759 Organization of the Primate Face Motor Cortex as Revealed by Intracortical. *J*
760 *Neurophysiol* 59:796–818.
- 761 Johnson K, Ladefoged P, Lindau M (1993) Individual differences in vowel production. *J*
762 *Acoust Soc Am* 94:701–714.
- 763 Kellis S, Miller K, Thomson K, Brown R, House P, Greger B (2010) Decoding spoken
764 words using local field potentials recorded from the cortical surface. *J Neural Eng*
765 7:56007–56016.
- 766 Kingyon J, Behroozmand R, Kelley R, Oya H, Kawasaki H, Narayanan NS, Greenlee
767 JDW (2015) High-gamma band fronto-temporal coherence as a measure of
768 functional connectivity in speech motor control. *Neuroscience* 305:15–25.
- 769 Kumar V, Croxson PL, Simonyan K (2016) Structural Organization of the Laryngeal
770 Motor Cortical Network and Its Implication for Evolution of Speech Production. *J*
771 *Neurosci* 36:4170–4181.
- 772 Ladefoged P, Johnson K (2011) *A Course in Phonetics*. Boston: Cengage Learning.
- 773 Lee DD, Seung HS (1999) Learning the parts of objects by non-negative matrix
774 factorization. *Nature* 401:788–791.
- 775 Leonard MK, Cai R, Babiak MC, Ren A, Chang EF (2016) Brain & Language The peri-
776 Sylvian cortical network underlying single word repetition revealed by
777 electrocortical stimulation and direct neural recordings. *Brain Lang*:1–15.
- 778 Leuthardt EC, Gaona C, Sharma M, Szrama N, Roland J, Freudenberg Z, Solis J,
779 Breshears J, Schalk G (2011) Using the electrocorticographic speech network to
780 control a brain-computer interface in humans. *J Neural Eng* 8.
- 781 Li M, Kambhmettu C, Stone M (2004) Automatic contour tracking in ultrasound
782 images. *Clin Linguist Phon* 19:545–554.
- 783 Li W, Chen Z, Yan N, Jones JA, Guo Z, Huang X (2016) Temporal Lobe Epilepsy Alters
784 Auditory-motor Integration For Voice Control. *Sci Rep* 6:1–13.
- 785 Lofqvist A, Gracco V (1999) Interarticulator programming in VCV sequences: Lip and
786 tongue movements. *J Acoust Soc Am* 105:1864–1876.
- 787 Lotze M, Seggewies G, Erb M, Grodd W, Birbaumer N (2000) The representation of
788 articulation in the primary sensorimotor cortex. *Neuroreport* 11:2985–2989.
- 789 Maeda S (1990) Compensatory articulation during speech: Evidence from the analysis
790 and synthesis of vocal-tract shapes using an articulatory model. In: *Speech*
791 *production and speech modeling*, pp 131–149. Netherlands: Springer.
- 792 Moran DW, Schwartz AB (1999) Motor Cortical Representation of Speed and Direction
793 During Reaching. *J Neurophysiol* 82:2676–2692.
- 794 Mugler EM, Patton JL, Flint RD, Wright Z a, Schuele SU, Rosenow J, Shih JJ,
795 Krusienski DJ, Slutzky MW (2014) Direct classification of all American English
796 phonemes using signals from functional speech motor cortex. *J Neural Eng*
797 11:35015.
- 798 Muller L, Hamilton LS, Edwards E, Bouchard KE, Chang EF (2016) Spatial resolution

- 799 dependence on spectral frequency in human speech cortex electrocorticography. *J*
800 *Neural Eng* 13.
- 801 Paninski L, Fellows MR, Hatsopoulos NG, Donoghue JP (2004) Spatiotemporal tuning of
802 motor cortical neurons for hand position and velocity. *J Neurophysiol* 91:515–532.
- 803 Pei X, Barbour DL, Leuthardt EC, Schalk G (2011) Decoding vowels and consonants in
804 spoken and imagined words using electrocorticographic signals in humans. *J Neural*
805 *Eng* 8.
- 806 Penfield W, Boldrey E (1937) Somatic motor and sensory representation in the cerebral
807 cortex of man as studied by electrical stimulation. *Brain* 60:389–443.
- 808 Penfield W, Roberts L (1959) *Speech and brain mechanisms*. Princeton.
- 809 Perkell JS, Matthies ML, Svirsky M a, Jordan MI (1993) Trading relations between
810 tongue-body raising and lip rounding in production of the vowel /u/: a pilot “motor
811 equivalence” study. *J Acoust Soc Am* 93:2948–2961.
- 812 Perkell JS, Nelson WL (1985) Variability in production of the vowels /i/ and /a/. *Acoust*
813 *Soc Am* 77:1889–1895.
- 814 Petersen S, Fox PT, Posner M, M M, Raichle M (1988) Positron emission tomographic
815 studies of the cortical anatomy of single-word Processing.
- 816 Ray S, Maunsell JHR (2011) Different origins of gamma rhythm and high-gamma
817 activity in macaque visual cortex. *PLoS Biol* 9.
- 818 Shadmehr R (2017) Distinct neural circuits for control of movement vs . holding still. *J*
819 *Neurophysiol* 117:1431–1460.
- 820 Shalit U, Zinger N, Joshua M, Prut Y (2012) Descending Systems Translate Transient
821 Cortical Commands into a Sustained Muscle Activation Signal. *Cereb Cortex*
822 22:1904–1914.
- 823 Simonyan XK, Ackermann H, Chang EF, Greenlee JD (2016) New Developments in
824 Understanding the Complexity of Human Speech Production. *J Neurosci* 36:11440–
825 11448.
- 826 Tourville J a, Guenther FH (2011) The DIVA model: A neural theory of speech
827 acquisition and production. *Lang Cogn Process* 26:952–981.
- 828 Ueda Y, Hamakawa T, Sakata T, Hario S, Watanabe A (2007) A real-time formant
829 tracker based on the inverse filter control method. *Acoust Sci Technol* 28:271–274.
- 830 Watanabe a. (2001) Formant estimation method using inverse-filter control. *IEEE Trans*
831 *Speech Audio Process* 9:317–326.

832
833
834
835

836 **Figure Legends**

837

838 **Figure 1: Experimental setup and articulator monitoring.** **A**, Schematic of the
839 articulatory tracking system. A video camera placed in front of the subject recorded the
840 movements of the lips while an ultrasound transducer under the jaw captured the tongue
841 contour. **B**, Example images of the from the video (top) and ultrasound (bottom) imaging
842 during production of the corner vowels /a/, /i/, and /u/. The lips and tongue contour were
843 extracted from these images, and the resulting binary masks are shown in color on top of
844 the raw images. **C**, Magnetic resonance imaging (MRI) reconstruction of the brains of

845 the four subjects included in the study. Co-registered ECoG electrodes are plotted on the
846 cortical surface, with dark points denoting electrodes over vSMC.

847

848 **Figure 2: Articulatory and acoustic behavior correlates with single electrode vSMC**
849 **neural activity.** **A**, Prototypical articulator positions for the corner vowels /a/, /i/, and
850 /u/. **B**, Average (\pm s.e.m.) timecourses of measured articulator displacements during
851 production of the corner vowels. For illustration, the kinematic parameter of vertical
852 position is shown, however both vertical and horizontal measurements are used for
853 subsequent analyses. **C**, Average formant values (F1-F4) for the corner vowels during the
854 same productions as in (B). **D**, HG activity in three example vSMC electrodes. Each
855 electrode was selected to be representative for the articulator shown in the top subplot
856 (median configurations shown for back tongue, front tongue, and lips). Trials for each
857 electrode are binned by the displacement of the articulator which best correlates with the
858 HG values at acoustic onset (\pm 100 ms). Yellow bars mark timepoints of significant
859 difference between the bins (e1: $F_{(2,75)} < 5.1$, $p < 0.01$; e2: $F_{(2,102)} < 4.9$, $p < 0.01$; e3:
860 $F_{(2,105)} < 4.8$, $p < 0.01$; ANOVA). **E** HG activity in the same electrodes binned according
861 to formant values (F1/F2 ratio).

862

863 **Figure 3: vSMC activity primarily encodes speech articulators.** **A**, Performance of
864 encoding models predicting vSMC HG using vowel identify, acoustic formants,
865 articulator kinematics, or all three. Articulator kinematics explain vSMC activity better
866 than vowel identity and acoustic formants (*** $p < 0.01$; Wilcoxon rank sum).
867 Furthermore, the joint model does not explain more variance than the kinematic model
868 alone, indicating that the vowel identity and acoustic formant models are likely driven by
869 variance shared with the kinematics. **B**, Electrodes over vSMC from three right
870 hemisphere subjects were warped onto a common brain and color-coded according to
871 articulator selectivity. Empty circles mark electrodes with no significant selectivity for
872 any articulator; black electrodes are selective for more than one articulator, and blue, red,
873 and green electrodes are selective for front tongue, back tongue, or lips, respectively.

874

875 **Figure 4: vSMC activity reflects articulator kinematic organization of vowels.** **A-C**,
876 Multidimensional Scaling (MDS) representations of **(A)** acoustic formants, **(B)** articulator
877 position, and **(C)** vSMC HG activity. Each letter marks the position of the median
878 production of the indicated vowel in MDS space across all subjects. The relative
879 organization of the vowels is similar across spaces. For example, the low-back vowel /a/
880 is always near the mid-back vowel / Λ /, but far from the high-front vowel /i/. **D-F**,
881 Hierarchical clustering (left) and confusion matrices (right) derived from the pairwise
882 distances between vowels in the MDS spaces of **(D)** acoustic formants, **(E)** articulator
883 position, and **(F)** vSMC HG.

884

885 **Figure 5: Representations of position, speed, velocity, and acceleration kinematic**
886 **features over time.** **A-D**, Top: Example kinematic parameters of position **(A)**, speed
887 **(B)**, velocity **(C)**, and acceleration **(D)** for all utterances of /a/ for one subject. Thin lines
888 mark individual trials, while the thick line is the across-trial average. Bottom:
889 Performance of encoding models predicting vSMC HG from articulator position **(A)**,
890 speed **(B)**, velocity **(C)**, or acceleration **(D)**. vSMC electrodes with significant

891 performance are marked by red dashes at the time peak encoding performance.
892 Electrodes in A-D are plotted in order of their peak encoding times in the joint model (F).
893 Vertical black lines mark the onset (solid) and offset (dashed) of vowel acoustics. **E**,
894 Comparison of the number of significant electrodes (black) and average peak
895 performance (grey) of position, speed, velocity, and acceleration encoding models. Speed
896 is significantly encoded at more electrodes, with a higher average model performance
897 (***p*<0.01; Wilcoxon rank sum). **F**, Performance of encoding models predicting vSMC
898 HG from all articulator kinematics jointly. **G**, Average performance across significant
899 electrodes for the joint and independent parameter models.

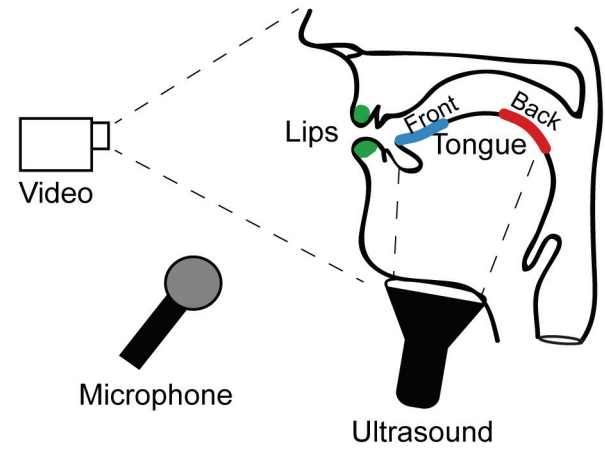
900
901 **Figure 6: Time course of kinematic parameter decoding.** **A-D**, Performance of
902 decoding models predicting articulator position (**A**), speed (**B**), velocity (**C**), or
903 acceleration (**D**) from HG at all vSMC electrodes. Features are averaged across subjects
904 and within articulators. Black lines denote the onset and offset of vowel acoustics.

905
906 **Figure 7: Relationship of vSMC HG dynamics and kinematic encoding.** **A**, HG
907 activity at several example electrodes illustrating the diverse dynamics during the same
908 behavior, especially during the time period when the vowel is being held. **B**, The first
909 two NMF bases extracted from HG dynamics across all electrodes. These bases
910 recapitulate the key differences in dynamics seen in the example electrodes, and serve as
911 an unbiased quantification of the HG dynamics seen across vSMC. **C**, HG activity at all
912 vSMC electrodes ordered by the ratio of NMF bases used to reconstruct their activity
913 ($(NMF1-NMF2)/(NMF1+NMF2)$). Arrow denotes the example electrode used in Figure
914 8. **D**, Average HG activity across all sustained ($NMF1>NMF2$) electrodes (solid) plotted
915 alongside the average encoding performance (dashed) across time. During the steady
916 state of the vowel, there is elevated activity, but almost no significant encoding of
917 articulator kinematics.

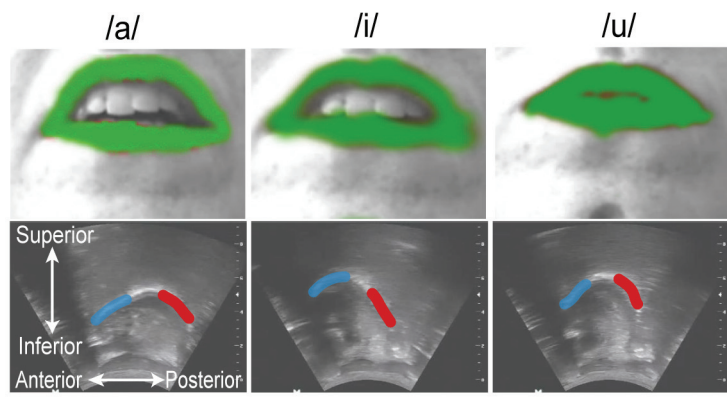
918
919 **Figure 8: Sustained HG activity is related to vowel duration.** **A**, An example
920 electrode that is characterized by sustained HG activity. Trials are ordered by vowel
921 production duration. Red lines mark the onset and offset of vowel acoustics. **B**, The
922 duration of the sustained HG activity at each trial is plotted against the trial duration for
923 all sustained electrodes. Larger grey markers denote observations from the example
924 electrode in A.

925
926

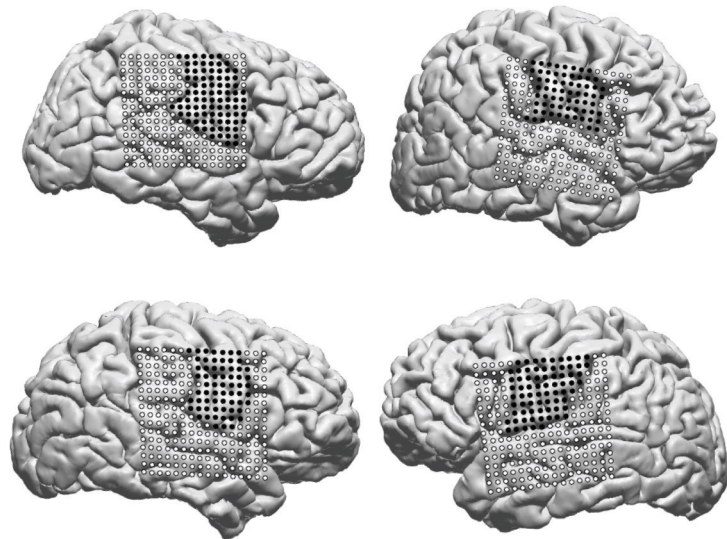
A Articulator Tracking Setup



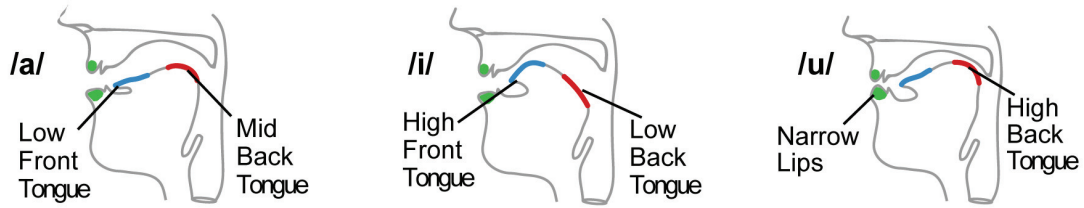
B Extracted Articulators



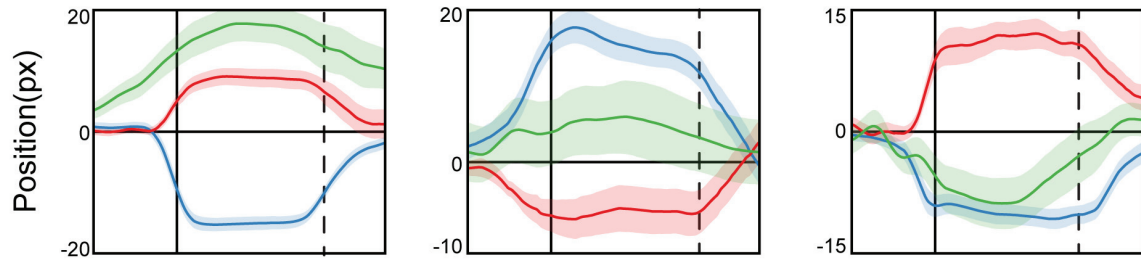
C ECoG Electrodes



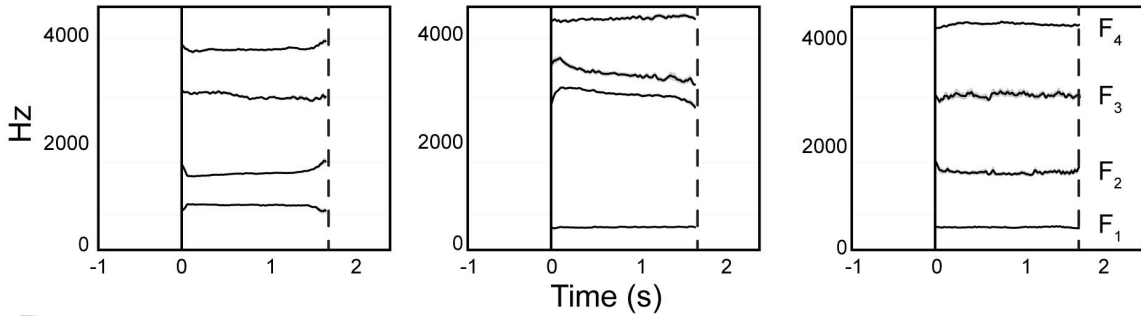
A Articulator Schematics



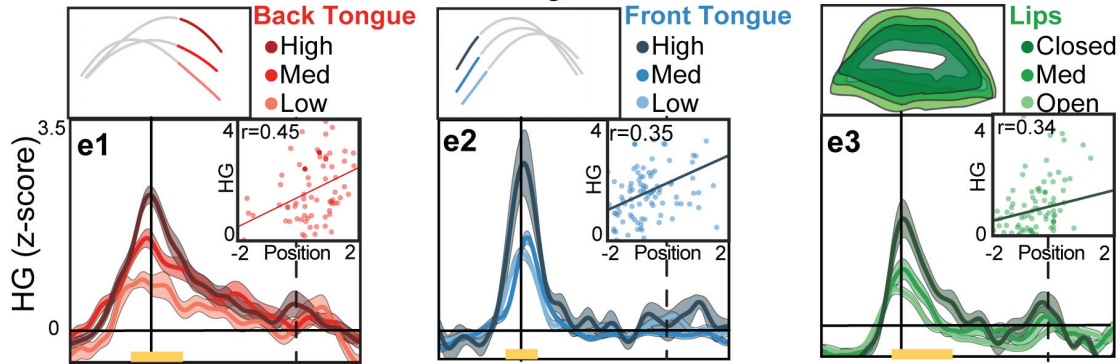
B Articulator Kinematics



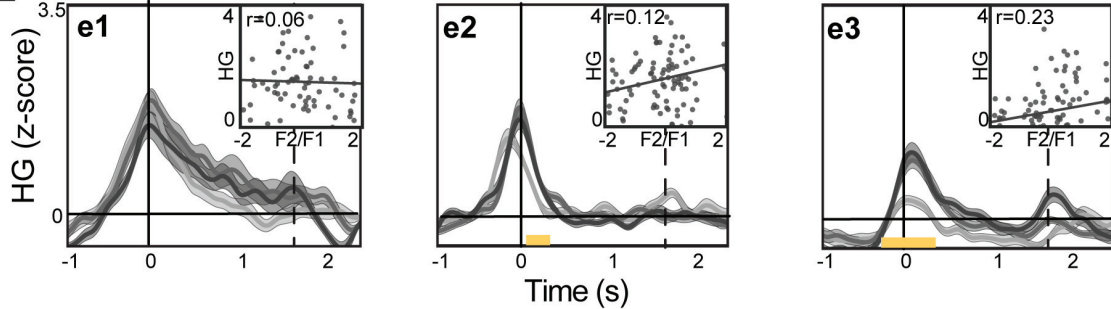
C Acoustic Formants



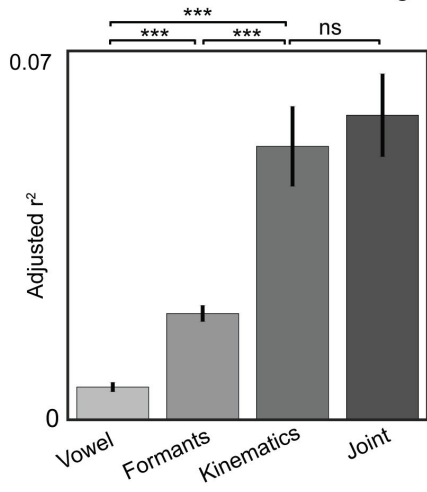
D vSMC electrode articulator encoding



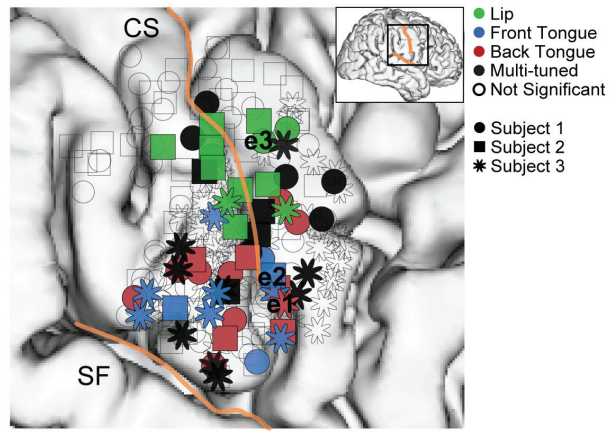
E vSMC electrode formant encoding

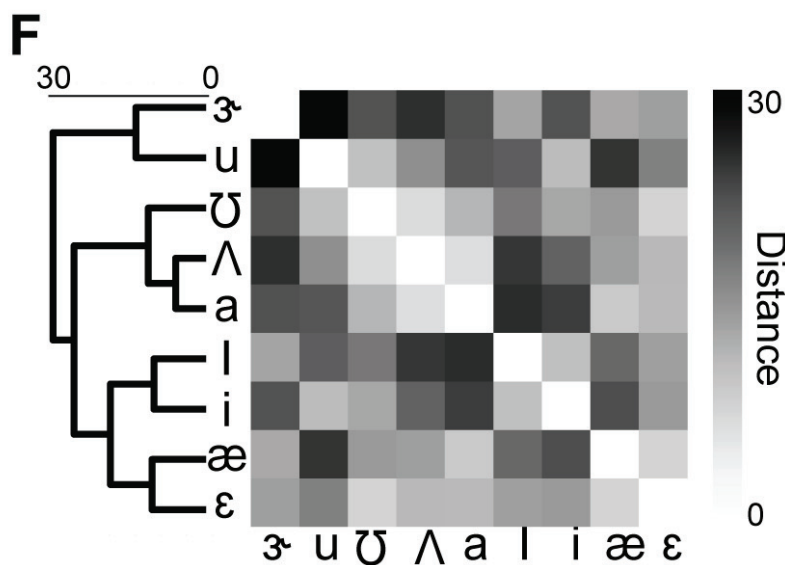
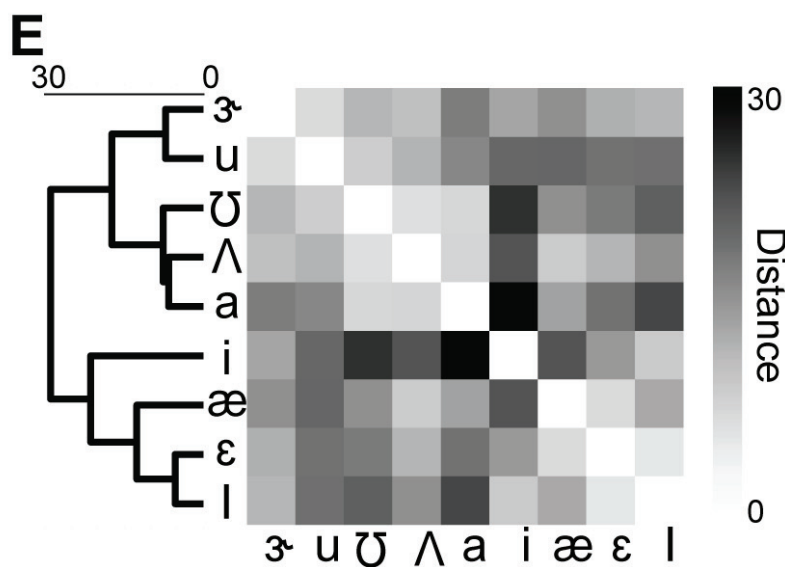
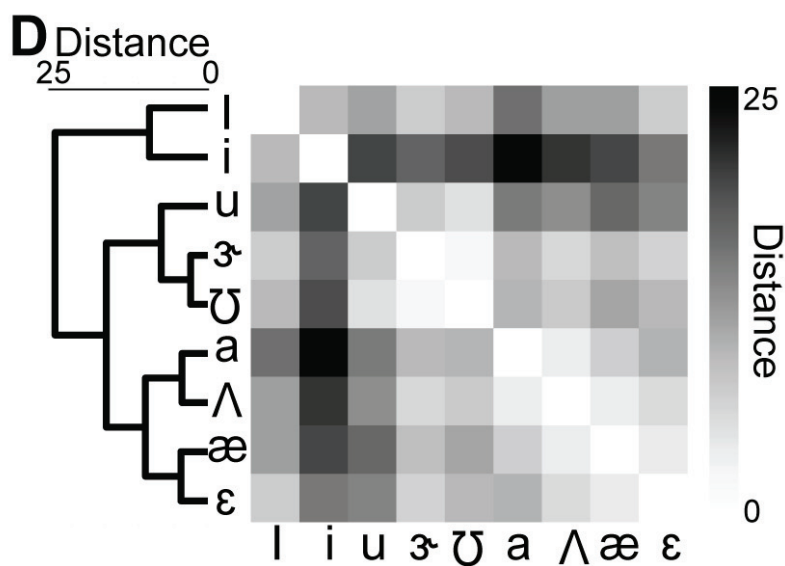
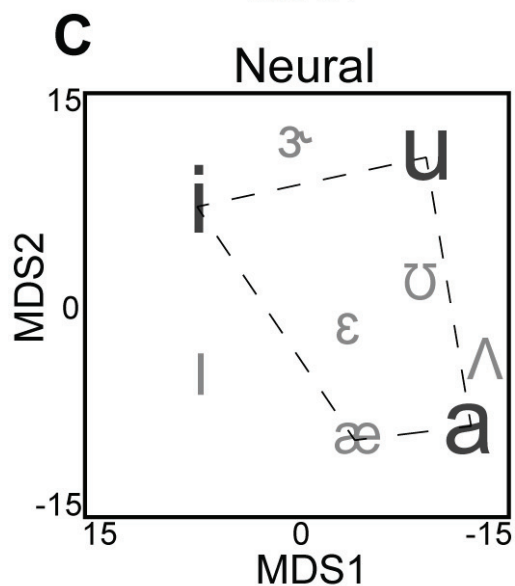
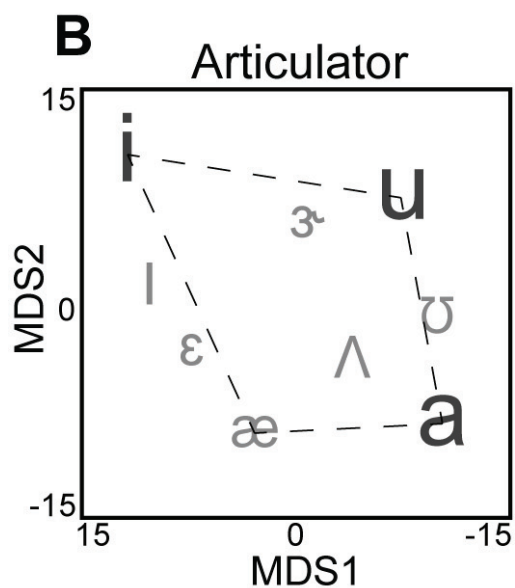
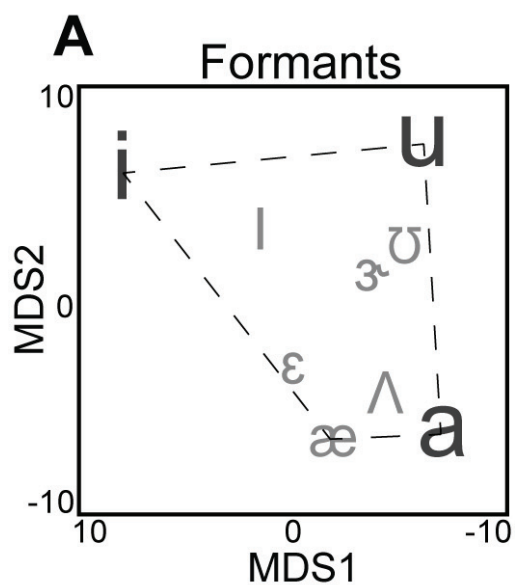


A Individual Electrode Encoding

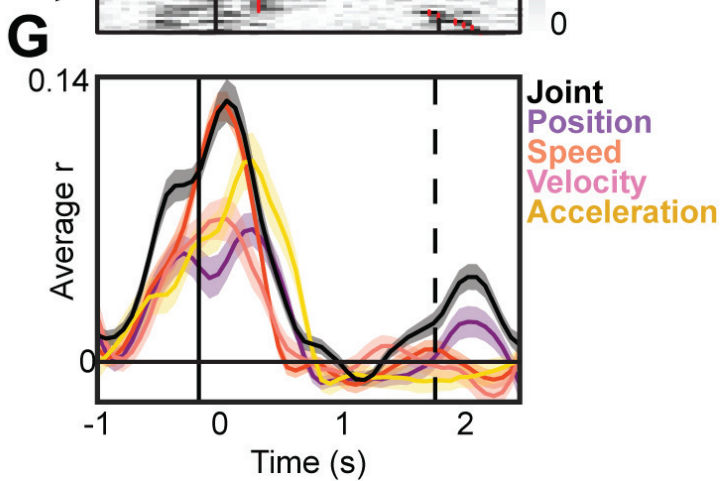
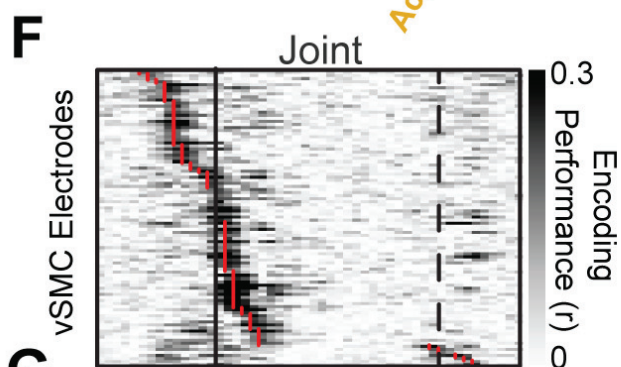
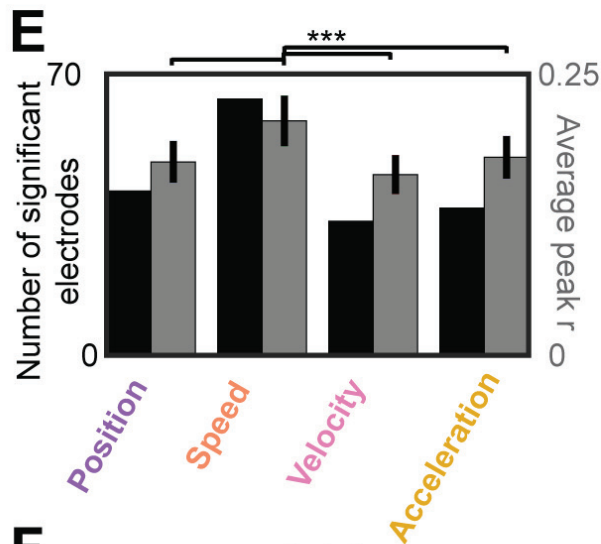
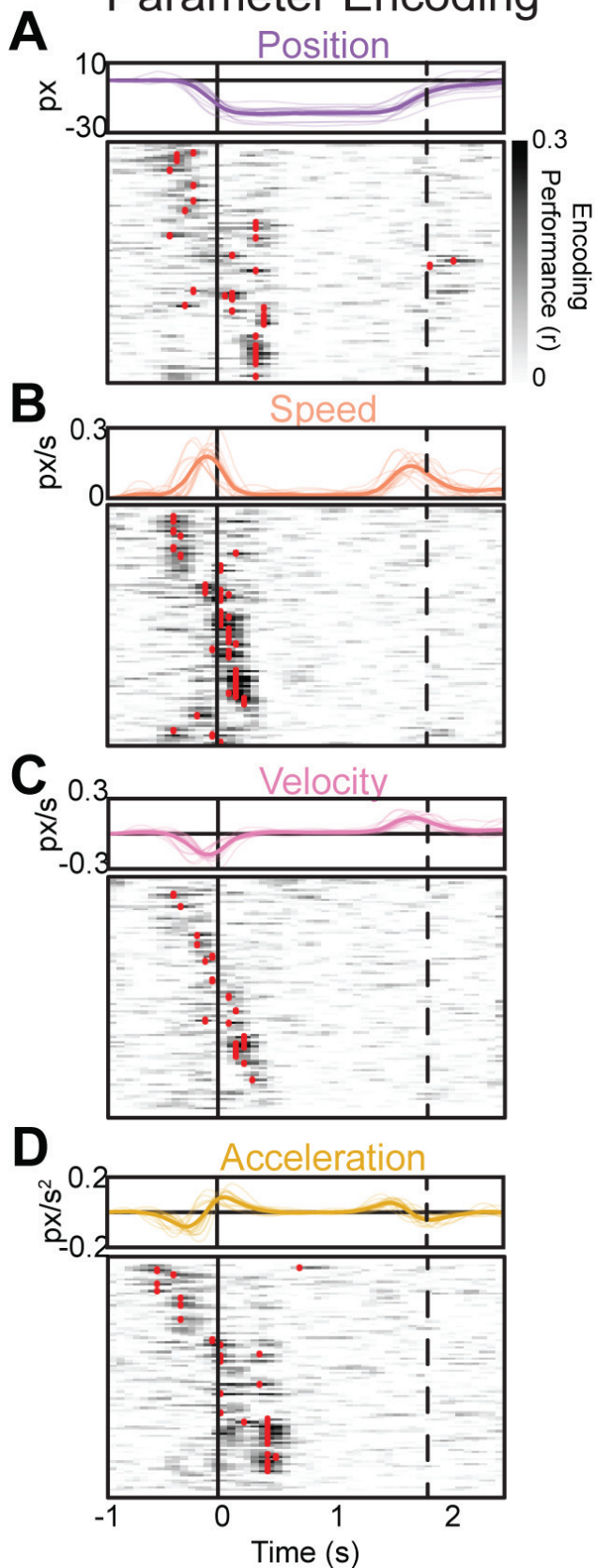


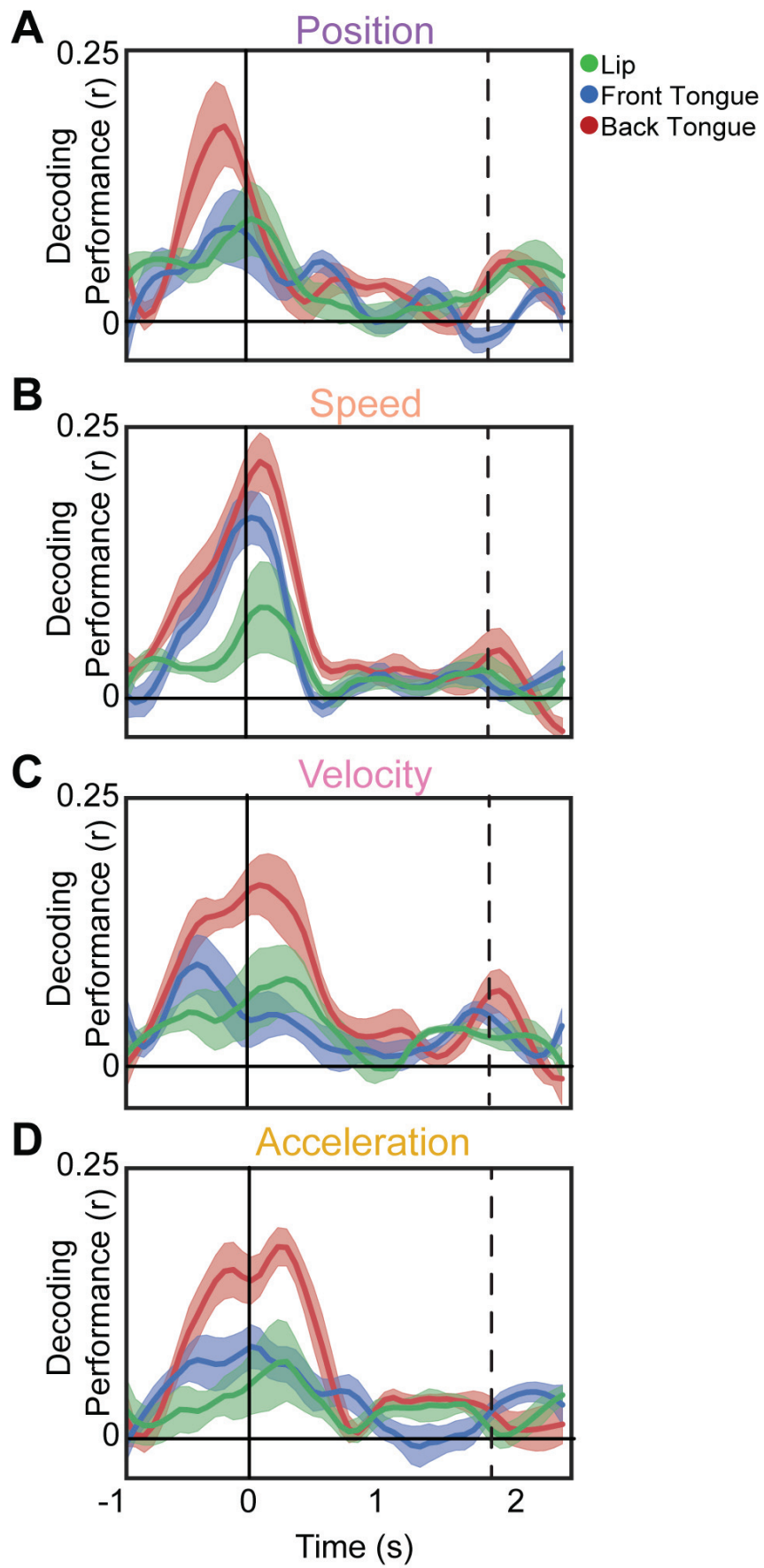
B vSMC Somatotopy



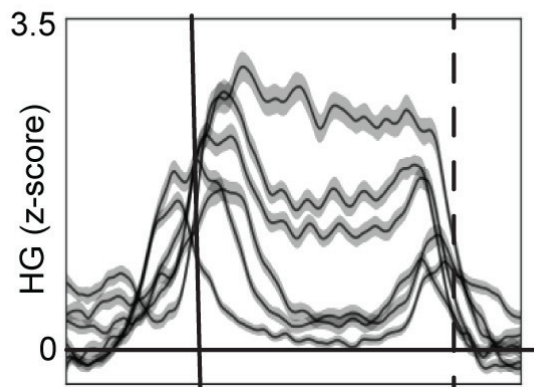


Timecourse of Kinematic Parameter Encoding

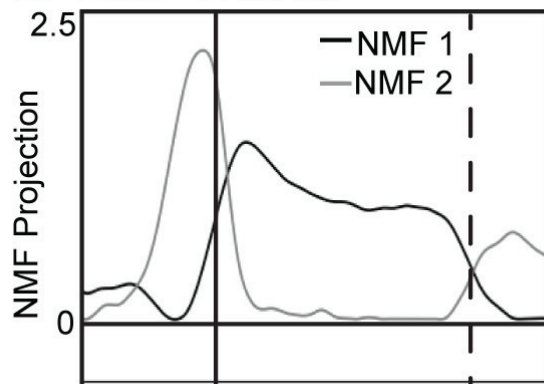




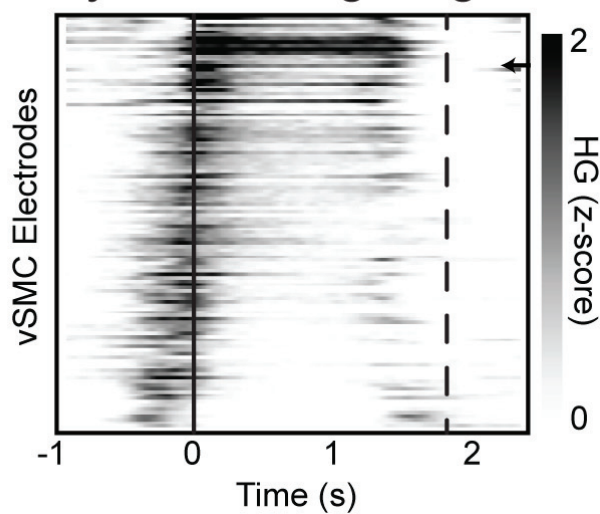
A Electrode Timecourses



B NMF Bases



C Electrodes Ordered by NMF Weighting



D Average Sustained HG vs Average Encoding

