

# 1 **Neural hierarchy for coding articulatory dynamics in speech imagery** 2 **and production**

3

4 Zehao Zhao<sup>1,2,3,4,5,#</sup>, Zhenjie Wang<sup>6,7,#</sup>, Yan Liu<sup>1,3,4,8,#</sup>, Youkun Qian<sup>1,3,4</sup>, Yuan Yin<sup>6</sup>, Xiaowei  
5 Gao<sup>9</sup>, Binke Yuan<sup>9</sup>, Shelley Xiuli Tong<sup>10</sup>, Xing Tian<sup>11,12,13</sup>, Gao Chen<sup>2,5,\*</sup>, Yuanning Li<sup>6,7,14,15,\*</sup>,  
6 Junfeng Lu<sup>1,3,4,16,\*</sup>, Jinsong Wu<sup>1,3,4,\*</sup>.

7

## 8 **Author affiliations:**

9 1 Department of Neurosurgery, Huashan Hospital, Shanghai Medical College, Fudan  
10 University, Shanghai 200040, China

11 2 Department of Neurosurgery, Second Affiliated Hospital, School of Medicine, Zhejiang  
12 University, Hangzhou 310016, China

13 3 Shanghai Key Laboratory of Clinical and Translational Brain-Computer Interface Research,  
14 Shanghai 200040, China

15 4 National Center for Neurological Disorders, Huashan Hospital, Shanghai Medical College,  
16 Fudan University, Shanghai 200040, China

17 5 Key Laboratory of Precise Treatment and Clinical Translational Research of Neurological  
18 Diseases of Zhejiang Province, Hangzhou 310016, China

19 6 School of Biomedical Engineering, ShanghaiTech University, Shanghai 201210, China

20 7 State Key Laboratory of Advanced Medical Materials and Devices, ShanghaiTech University,  
21 Shanghai 201210, China

22 8 Department of Neurological Diagnosis and Restoration, Osaka University Graduate School  
23 of Medicine, Suita, Osaka 565-0871, Japan

24 9 Institute for Brain Research and Rehabilitation, South China Normal University, Guangzhou  
25 510631, China

26 10 Human Communication, Learning, and Development, Faculty of Education, The University  
27 of Hong Kong, Hong Kong 999077, China

28 11 New York University Shanghai, Shanghai 200031, China

29 12 Shanghai Key Laboratory of Brain Functional Genomics (Ministry of Education), School  
30 of Psychology and Cognitive Science, East China Normal University, Shanghai 200031, China

31 13 New York University-East China Normal University Institute of Brain and Cognitive  
32 Science at New York University Shanghai, Shanghai 200031, China.

33 14 Shanghai Clinical Research and Trial Center, Shanghai 201210, China

34 15 Lingang Laboratory, Shanghai 200031, China.

35 16 MOE Frontiers Center for Brain Science, Huashan Hospital, Fudan University, Shanghai  
36 200040, China.

37 #These authors contributed equally to this work.

38 \*Correspondence and requests for materials should be addressed to: Jinsong Wu  
39 ([wujinsong@huashan.org.cn](mailto:wujinsong@huashan.org.cn)), Junfeng Lu ([junfeng\\_lu@fudan.edu.cn](mailto:junfeng_lu@fudan.edu.cn)), Yuanning Li  
40 ([liyn2@shanghaitech.edu.cn](mailto:liyn2@shanghaitech.edu.cn)), or Gao Chen ([d-chengao@zju.edu.cn](mailto:d-chengao@zju.edu.cn)).

41 **Abstract**

42 Mental imagery is a hallmark of human cognition, yet the neural mechanisms underlying these  
43 internal states remain poorly understood. Speech imagery—the internal simulation of speech  
44 without overt articulation—has been proposed to partially share neural substrates with actual  
45 speech articulation. However, the precise feature encoding and spatiotemporal dynamics of this  
46 neural architecture remain controversial, constraining the understanding of mental states and  
47 the development of reliable speech imagery decoders. Here, we leveraged high-resolution  
48 electrocorticography recordings to investigate the shared and modality-specific cortical coding  
49 of articulatory kinematic trajectories (AKTs) during speech imagery and articulation. Applying  
50 a linear model, we identified robust neural dynamics in frontoparietal cortex that encoded AKTs  
51 across both modalities. Shared neural populations across the middle premotor cortex,  
52 subcentral gyrus, and postcentral-supramarginal junction exhibited consistent spatiotemporal  
53 stability during the integrative articulatory planning. In contrast, modality-specific populations  
54 for speech imagery and articulation were somatotopically interleaved along the primary  
55 sensorimotor cortex, revealing a hierarchical spatiotemporal organization distinct from shared  
56 encoding regions. We further developed a generalized neural network to decode multi-  
57 population neural dynamics. The model achieved high syllable prediction accuracy for speech  
58 imagery (79% median accuracy), closely matching the performance of speech articulation  
59 (81%). This model robustly extrapolated AKT decoding to untrained syllables within each  
60 modality while demonstrating cross-modal generalization across shared populations. These  
61 findings uncover a somato-cognitive hierarchy linking high-level supramodal planning with  
62 modality-specific neural manifestation, revolutionizing an imagery-based brain-computer  
63 interface that directly decodes thoughts for synthetic telepathy.

64

65 **Keywords:** speech imagery, speech articulation, articulatory kinematic trajectories, high-  
66 density electrocorticography, frontoparietal cortex, neural coding

67

## 68 **Introduction**

69 Mental imagery is a unique adaptive trait of human cognition for linking past experiences,  
70 current states, and future scenarios by simulating and creating mental events<sup>1-5</sup>. For example,  
71 mental imagery of speech, also known as covert or inner speech, represents the internally  
72 generated, quasi-perceptual experience of speech without overt or audible articulation<sup>6,7</sup>. This  
73 phenomenon exemplifies both an instantiation of language as a vehicle for thought and a  
74 specific form of environment-disengaged action<sup>8</sup>. Existing research has shown that the neural  
75 processes of overt speech articulation involve a multi-stage sequence that comprises conceptual  
76 formulation, lexical selection, phonological encoding, motor planning, and the execution of  
77 articulatory movements<sup>6,7,9,10</sup>. However, the neural mechanisms governing speech imagery are  
78 not yet fully understood. Current theories propose that speech imagery may constitute a  
79 truncated form of overt speech production, but the specific stage at which this process diverges  
80 remains debated<sup>6,11-13</sup>. The abstraction view posits that speech imagery arises from the  
81 activation of abstract linguistic constructs, independent of articulatory representations<sup>6,11,14,15</sup>.  
82 In contrast, the motor simulation hypothesis suggests that speech imagery mirrors overt speech  
83 articulation, involving the planning and prediction of articulatory details that are halted before  
84 execution<sup>6,11,13</sup>. Clarifying this debate is critical for understanding how speech motor control  
85 interfaces with internal cognitive states.

86 Emerging neurolinguistic models, supported by findings from functional magnetic  
87 resonance imaging (fMRI), magnetoencephalography (MEG), and scalp  
88 electroencephalography (EEG), suggest that initial processes involved in motor planning for  
89 both speech imagery and articulation exhibit considerable similarities<sup>7,16-18</sup>. Functional  
90 localization studies reveal overlapping frontoparietal activation in regions such as the premotor  
91 cortex and supramarginal gyrus (SMG), along with the articulation-specific engagement of the  
92 primary sensorimotor cortex (SMC)<sup>6,16-21</sup>. Compared to actual articulation, the neural  
93 operations during speech imagery appear to suspend downstream output execution, retaining  
94 only an efference copy of motor commands<sup>16-18</sup>. This copy is passed to an internal forward  
95 model in the frontoparietal cortex to facilitate internal sensorimotor simulation and predict the

96 subsequent outcomes of speech<sup>16-18</sup>. However, the precise neural hierarchy encoding internal  
97 and external articulatory movements remains unsolved, primarily due to the limited  
98 spatiotemporal resolution of non-invasive techniques and a paucity of high-precision  
99 intracranial recordings during speech imagery.

100 Beyond theoretical implications, speech imagery offers potential for neural decoding,  
101 particularly in restoring communication for individuals with severe speech impairments (e.g.,  
102 stroke, amyotrophic lateral sclerosis [ALS], or locked-in syndrome)<sup>22</sup>. Recent intracranial  
103 EEG studies have made strides in decoding vocalized, attempted, and mimed speech by  
104 leveraging articulatory movement encoding in the frontoparietal cortex<sup>23-31</sup>. Yet, these  
105 approaches rely on residual motor execution signals, limiting their utility for patients with  
106 complete loss of speech output<sup>22,32</sup>. Decoding speech imagery may overcome this constraint  
107 by targeting internally generated speech processes, presenting a more adaptable solution for  
108 communication restoration.

109 Here, we employed high-density electrocorticography (ECoG) recordings from nine  
110 participants to dissect the encoding of articulatory kinematic trajectories (AKTs) during  
111 syllable articulation (SA) and syllable imagery (SI). AKTs transcend categorical approaches by  
112 providing dynamic articulatory representations that both resolve theoretical debates  
113 (abstraction vs. motor simulation) and establish a biologically grounded foundation for  
114 naturalistic speech decoding<sup>28</sup>. Using a temporal receptive field (TRF) model<sup>28,33,34</sup>, we  
115 identify modality-specific AKT encoding preferences across frequency bands, characterize  
116 shared and distinct neural substrates, and map their spatiotemporal and somatotopic  
117 organization. Finally, we developed a unified decoding framework for both modalities' syllable  
118 classification, speech synthesis, and AKT reconstruction. Notably, this framework  
119 demonstrates robust cross-modal generalization, and successfully extrapolates AKT decoding  
120 to untrained syllables. Our findings reveal a hierarchical somato-cognitive architecture  
121 underpinning speech imagery and articulation while advancing clinically viable brain-  
122 computer interfaces (BCIs) that decode speech via mental imagery in individuals with a broad  
123 spectrum of speech output impairments.

## 124 **Results**

### 125 **Overview of experiments**

126 Nine subjects (S1-S9) were recruited in this study. All subjects received awake surgeries for  
127 the treatment of brain tumors. During the awake surgery, the subjects performed the syllable  
128 articulation (SA)—syllable imagery (SI) contrast task while simultaneous audio and 256-  
129 channel high-density ECoG recordings were obtained. Among the subjects, five (S1-S5)  
130 completed Paradigm 1 (**Fig. 1a**), which involved 60 repetitions of articulation and silent  
131 imagery of six syllables, each initiated by an ipsilateral hand key press. To control for any  
132 potential influence of hand motor responses, 40 repetitions of ipsilateral key presses were  
133 performed. Electrodes responsive solely to these key presses were excluded from further  
134 analyses. The remaining four subjects (S6-S9) completed Paradigm 2 (**Fig. 1b**), which  
135 consisted of 45 repetitions of articulation and imagery tasks across eight syllables, each  
136 initiated by the appearance of a black cross symbol. To ensure the complete absence of motor  
137 and auditory output during speech imagery tasks, we implemented preoperative training  
138 combined with continuous intraoperative video/audio surveillance and electromyographic  
139 (EMG) monitoring. (Detailed in the **Methods** section)

140

### 141 **Distinct spectral signatures govern AKT encoding across SA and SI modalities**

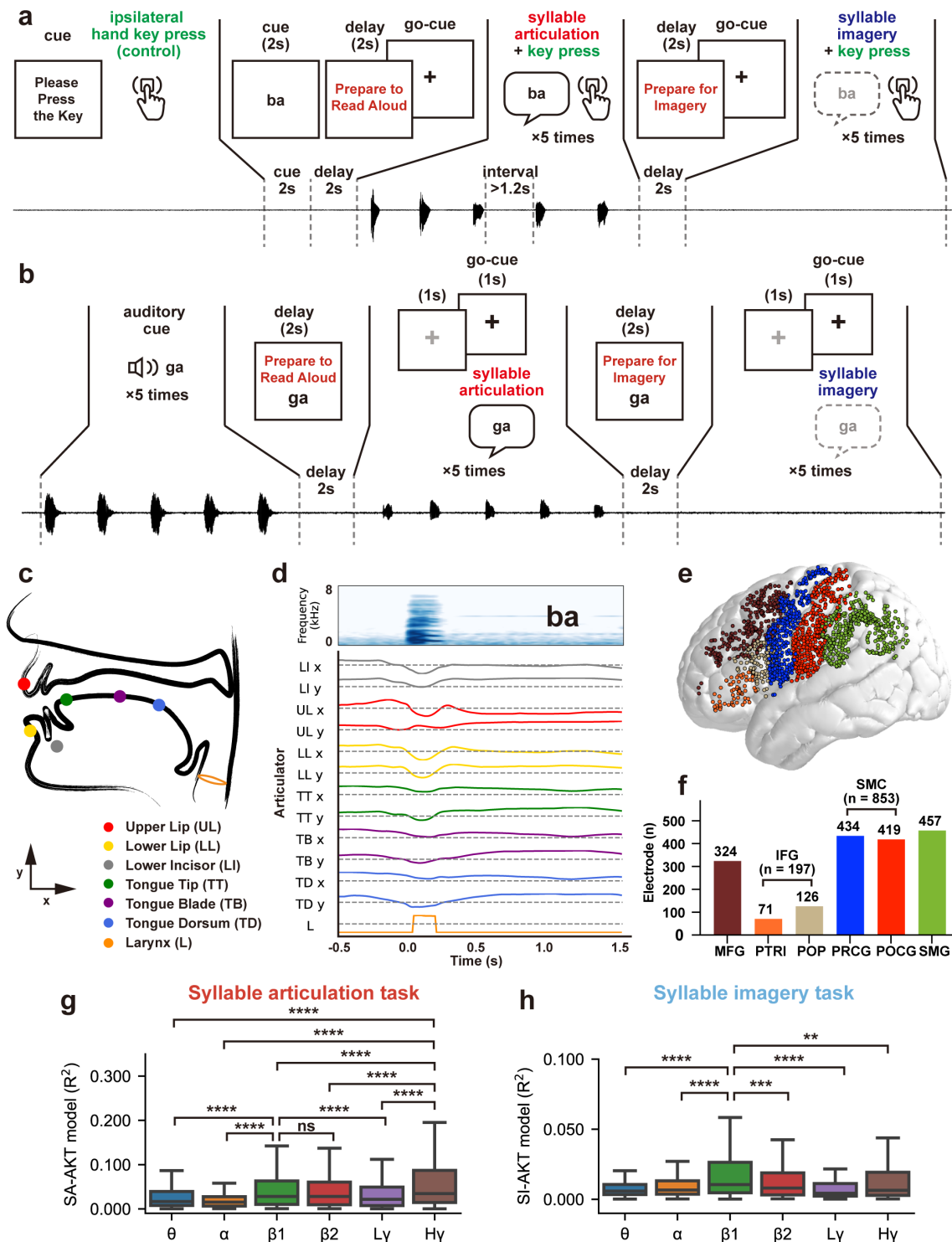
142 We first analyzed the neural representation of the AKTs during the SA and SI tasks. Using an  
143 established speaker-independent acoustic-articulatory inversion (AAI) algorithm<sup>28</sup>, we derived  
144 13-dimensional AKTs from synchronized audio recordings during SA (**Fig. 1c-d**). These  
145 trajectories captured: (i) two-dimensional movements (x, y) of six articulators (upper/lower  
146 lips, jaw, tongue tip/body/dorsum); and (ii) one-dimensional laryngeal activity (scaled log F<sub>0</sub>)  
147<sup>28</sup>. For SI tokens, we used the corresponding SA-derived AKTs as internal kinematic proxies.

148 Subsequently, we computed the analytic amplitude of the signals from the local field  
149 potential across seven frequency bands—theta (4-8 Hz), alpha (8-12 Hz), beta1 (12-24 Hz),  
150 beta2 (24-40 Hz), low gamma (40-70 Hz), and high gamma (70-150 Hz)<sup>21,25</sup>. This analysis  
151 encompassed 1,831 electrodes distributed across our anatomically defined regions of interest

152 (ROIs) (434 over precentral gyrus [PRCG], 419 over postcentral gyrus [POCG], 126 over pars  
153 opercularis [POP], 71 over pars triangularis [PTRI] of Broca's area, 324 over middle frontal  
154 gyrus [MFG], and 457 over SMG) in nine participants (**Fig. 1e-f**). The beta1 band showed  
155 maximal modality responsiveness (SA: 60.0%; SI: 40.6% of electrodes), exceeding other bands  
156 (**Extended Data Fig. 1**).

157 To quantify how AKTs are represented in cortical activity, we modeled the relationship  
158 between time-varying AKTs and neural signals using a time-delayed ridge regression (TRF  
159 model) <sup>28,33,34</sup>. This approach revealed striking modality-dependent encoding profiles: SA  
160 encoding peaked in high gamma band (median  $R^2$ : 0.034, interquartile range [IQR]: 0.014–  
161 0.087; all  $p < 0.0001$  vs other bands, Mann-Whitney  $U$  tests, Benjamini-Hochberg corrected),  
162 with secondary beta1 band representation (median  $R^2$ : 0.028, IQR: 0.010–0.063;  $p < 0.0001$   
163 vs theta, alpha, and low gamma bands,  $p = 0.3784$  vs beta2) (**Fig. 1g**). SI encoding was beta1-  
164 dominant (median  $R^2$ : 0.011, IQR: 0.005–0.026; all  $p < 0.01$  vs other bands, **Fig. 1h**). This  
165 spectral dissociation persisted across all frontoparietal regions (**Extended Data Fig. 2**),  
166 suggesting fundamental differences in how internal simulation versus motor execution engages  
167 the speech network.

168



169

170 **Fig. 1 | Spectral dissociation in articulatory kinematic trajectory (AKT) encoding between**  
 171 **syllable articulation (SA) and syllable imagery (SI).**

172 **a**, Schematic of Paradigm 1. The task begins with a “Press Key” prompt, requiring subjects to press  
 173 a key with the ipsilateral hand (green) 40 times as a baseline control. A visual syllable stimulus  
 174 (e.g., “ba”) is then presented for 2 s, followed by a 2-s “Prepare to Read Aloud” prompt. After the  
 175 appearance of the black cross symbol (go-cue), subjects articulate the syllable five times (red, SA

176 task), marking onset with an ipsilateral key press (green, interval > 1.2 s). Next, a 2-s “Prepare for  
177 Imagery” prompt is shown, followed by the black cross symbol (go-cue), after which subjects  
178 imagine the syllable five times (blue, SI task), again marking onset with an ipsilateral key press  
179 (green, interval > 1.2 s). **b**, Schematic of Paradigm 2. Subjects first listen to an auditory syllable  
180 stimulus (e.g., “ga”) repeated five times. A 2-s “Prepare to Read Aloud” prompt and visual cue are  
181 followed by alternating grey and black crosses (1 s each). Subjects articulate the syllable upon the  
182 black cross (red, SA task), repeated five times. For imagery, a 2-s “Prepare for Imagery” prompt  
183 and visual cue precede alternating grey and black crosses, with subjects imagining the syllable  
184 upon each black cross (blue, SI task), repeated five times. **c**, Diagram of AKT movements for six  
185 articulators: upper lip (UL), lower lip (LL), lower incisor/jaw (LI), tongue tip (TT), tongue blade (TB),  
186 and tongue dorsum (TD) in 2-dimensional plane (x and y axes), along with laryngeal movement (L,  
187  $\log F_0$ , scaled to -1 to 1 arbitrary units [a.u.],  $F_0$  refers to the fundamental frequency). **d**, Example  
188 Mel-spectrogram (upper) of “ba” converted into a 13-dimensional AKT time series (lower) via the  
189 acoustic-to-articulatory inversion (AAI) algorithm. **e**, The normalized total electrode coverage  
190 among cortical regions of interest on the MNI 152 brain template (subject  $N = 9$ , electrode  $n =$   
191 1,831): middle frontal gyrus (MFG, brown), pars triangularis (PTRI, orange), pars opercularis (POP,  
192 light yellow) of the inferior frontal gyrus (IFG), sensorimotor cortex (SMC), including precentral  
193 gyrus (PRCG, blue) and postcentral gyrus (POCG, red), and supramarginal gyrus (SMG, green).  
194 **f**, Bar graph depicting electrode counts per cortical region. **g-h**, Box plots of AKT model  
195 performance ( $R^2$ ) across frequency bands for SA (**g**) and SI (**h**): theta ( $\theta$ , 4-8 Hz), alpha ( $\alpha$ , 8-12  
196 Hz), beta1 ( $\beta_1$ , 12-24 Hz), beta2 ( $\beta_2$ , 24-40 Hz), low gamma (Ly, 40-70 Hz), and high gamma (Hy,  
197 70-150 Hz). \*\* $p < 0.01$ , \*\*\* $p < 0.001$ , \*\*\*\* $p < 0.0001$ , ns: non-significant (two-sided Mann-Whitney  
198  $U$  tests with Benjamini-Hochberg correction). Comparisons: SA (Hy or  $\beta_1$  vs. other bands) and SI  
199 ( $\beta_1$  vs. other bands).

200

201 **Frontoparietal beta1 activity underpins shared and specific AKT encoding in SA and SI**  
202 **modalities**

203 Given our findings on the frequency-band encoding preferences, along with extensive prior  
204 research on high gamma encoding in overt speech<sup>19-21,25,27,28</sup>, our analysis focuses on the beta1  
205 band to gain insights into speech imagery. We found that 56.9% of the 1287 responsive  
206 electrodes were responsive exclusively to SA (42.2%) or to SI (14.7%) in the beta1 band, and  
207 43.1% of all responsive electrodes responded to both modalities (**Fig. 2a**, see **Extended Data**  
208 **Fig. 3** for individual-level spatial distributions). Anatomical distributions differed significantly  
209 ( $p < 0.0001$ ,  $\chi^2$  test), with dual-responsive electrodes most prevalent in PRCG (47.9%),  
210 followed by POCG (33.9%) and SMG (22.8%) (**Fig. 2a**).

211 Applying stringent selection criteria ( $R^2 > 0.01$ , top 50% of responsive electrodes;  
212 **Extended Data Fig. 4b**), we identified functionally distinct neural populations exhibiting  
213 differential encoding of AKTs across speech modalities. SA-specific electrodes exhibited beta1  
214 activity that precisely tracked actual articulator movements. For example, Electrode E1 (**Fig.**  
215 **2b**, light red hexagon with black border) showed activity patterns (**Fig. 2g-h**) that scaled with  
216 tongue retraction demands, being strongest for alveolar /da/, intermediate for /du/, and weakest  
217 for bilabial /ba/ (**Fig. 2f**). These electrodes remained silent during imagery (**Fig. 2g-h**),  
218 revealing their exclusive role in motor execution.

219 Strikingly, we identified a distinct population of SI-specific encoding electrodes that  
220 maintained precise encoding of AKTs despite a complete absence of movement. Representative  
221 Electrode E2 (**Fig. 2b**, light blue hexagon with black border) exhibited significant beta1 activity  
222 (**Fig. 2j-k**) that tracked the internal dynamics of tongue tip elevation (**Fig. 2i**) during mental  
223 imagery. Specifically, the double-peaked beta1 activity pattern mirrored the kinematic  
224 sequence of simulated alveolar contact (/d/) followed by vowel-specific repositioning (**Fig. 2i-**  
225 **k**), demonstrating accurate internal simulation of articulation.

226 Our investigation of dual-responsive electrodes revealed two distinct patterns of AKT  
227 encoding across tasks. Quantitative clustering analysis (Duda-Hart statistic:  $d > 1.645$ ,  $p < 0.05$ ;  
228 silhouette score-optimized, **Extended Data Fig. 4a**) identified two functionally segregated

229 populations within the  $R^2$  value distribution space (SA-AKT vs SI-AKT models, **Fig. 2c**).  
230 The first subset of electrodes (Cluster 1, pink dots in **Fig. 2c-d**) exhibited a strong positive  
231 correlation between SA-AKT and SI-AKT  $R^2$  values (Pearson's  $r = 0.74$ ,  $p < 0.0001$ ),  
232 suggesting shared encoding across SA and SI modalities. Representative Electrode E3 (**Fig. 2b**,  
233 light purple hexagon with pink border) maintained consistent beta1 activity for sequential jaw  
234 elevation planning of /ba/, /da/, and /du/ syllables across both modalities ( $r = 0.95$  for /ba/,  $r =$   
235  $0.87$  for /da/,  $r = 0.68$  for /du/, all  $p < 0.0001$ , **Fig. 2n**). These populations faithfully tracked  
236 intended articulatory demands, showing enhanced activity for open vowel /a/ compared to close  
237 vowel /u/ (**Fig. 2l-m**).

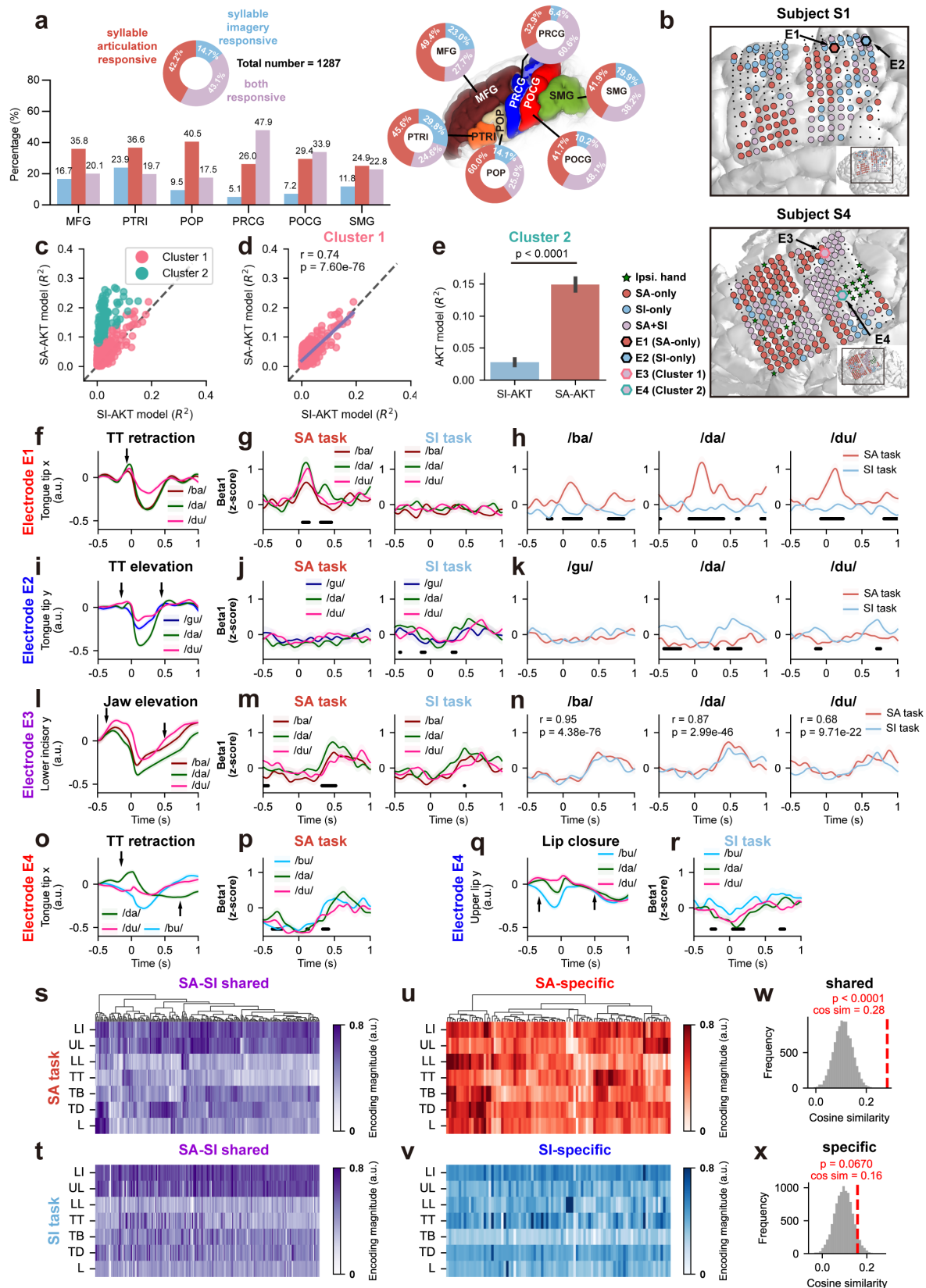
238 The other subset of electrodes (Cluster 2, light green dots in **Fig. 2c**), on the contrary,  
239 showed significant but distinct AKT encoding performances between SA and SI tasks ( $R^2$ :  
240  $0.149 \pm 0.005$  vs.  $0.028 \pm 0.002$ ,  $p < 0.0001$ , two-sided paired  $t$ -test, **Fig. 2e**). Representative  
241 Electrode E4 (**Fig. 2b**, light purple hexagon with light green border) switched functional  
242 specialization between modalities: during SA task it encoded tongue tip retraction (alveolar /d/ >  
243 bilabial /b/; **Fig. 2o-p**), while during SI task it preferentially represented lip closure (bilabial  
244 /b/ > alveolar /d/; **Fig. 2q-r**). This functional switching occurred with precise temporal coupling  
245 to expected articulator kinematics in each modality (**Fig. 2p, r**).

246 Taking a closer look at the coding properties in individual electrodes from these two  
247 clusters, we further confirmed these distinct SA-SI AKT coding patterns. For each cluster, we  
248 calculated each articulator's unique contribution ( $\Delta R^2$ ) to the full AKT model across electrodes.  
249 Cluster 1 encoding electrodes (e.g., Electrode E3) maintained robust cross-modal consistency  
250 between the SA modality (**Fig. 2s**) and the SI modality (**Fig. 2t**), with a *cosine similarity* of  
251  $0.28$  (**Fig. 2w**,  $p < 0.0001$ , permutation test). This stability suggests these neural ensembles  
252 participate in a supramodal speech planning network that transcends production modality.  
253 Conversely, in Cluster 2 encoding electrodes, no significant similarity was observed between  
254 the AKT encoding patterns of the SA modality (**Fig. 2u**) and the SI modality (**Fig. 2v**) (**Fig. 2x**,  
255 *cosine similarity* =  $0.16$ ,  $p = 0.0674$ , permutation test). This modality-specific specialization  
256 likely reflects microanatomical segregation of functionally distinct neural subpopulations  
257 within the same cortical territory, with differential engagement based on behavioral state (e.g.,

258 Electrode E4, serving as SA-specific electrodes in SA modality and as SI-specific electrodes  
259 in SI modality) <sup>35</sup>.

260 Our results demonstrate that speech imagery constitutes an active cognitive process  
261 involving detailed internal articulatory representations, distinct from truncated articulation.  
262 While partial overlap with speech articulation (Cluster 1, Electrode E3) confirms shared  
263 planning mechanisms <sup>6,7,9,10</sup>, the discovery of modality-specific neural signatures—particularly  
264 SI-exclusive AKT encoding (Electrode E2) and Cluster 2’s functional switching (Electrode  
265 E4)—reveals unique computational principles governing imagined speech. Identifying these  
266 internal kinematic codes provides a neurophysiological basis for decoding imagined speech,  
267 with immediate implications for brain-computer interfaces to restore communication in  
268 paralysis.

269



270

271 **Fig. 2 | Shared and specific AKT encoding in SA and SI modalities mediated by**  
 272 **frontoparietal beta1 activity.**

273 **a**, Electrode classification and anatomical distribution. Pie chart (top left) shows proportions of  
274 modality-selective (SA-only: red; SI-only: blue) and dual-responsive (purple) electrodes ( $n = 1,287$ ).  
275 Bar plot (bottom left) displays the regional distribution. Right pie charts depict gyrus-specific  
276 distributions. **b**, Individual brain reconstructions for subjects S1 and S4, illustrating electrode  
277 coverage. Green stars: ipsilateral hand motor-responsive electrodes; light red, light blue, and light  
278 purple circles: SA-only, SI-only, and dual-responsive electrodes, respectively. Hexagons mark  
279 representative electrodes: E1 (SA-only, light red), E2 (SI-only, light blue), E3 (dual-responsive,  
280 Cluster 1, light purple with pink outline), and E4 (dual-responsive, Cluster 2, light purple with green  
281 outline). Black circles: non-responsive electrodes. **c**, Scatter plot illustrates the relationship  
282 between the total variance explained ( $R^2$ ) of the SA-AKT and SI-AKT models for each dual-  
283 responsive electrode, with colors representing the k-means clustering classification (Cluster 1: pink;  
284 Cluster 2: light green). **d**, Scatter plot highlights a strong correlation between SA- and SI-AKT  
285 model performance for Cluster 1 electrodes ( $r = 0.74$ ,  $p = 7.60 \times 10^{-78}$ ), with a purple fitted curve  
286 closely aligned to the  $y = x$  diagonal (black dashed line). **e**, Bar plot shows significant differences  
287 in SA- and SI-AKT model performance ( $R^2$ , mean  $\pm$  s.e.m.) for Cluster 2 electrodes ( $p < 0.0001$ ,  
288 two-sided paired  $t$ -test). **f-h**, Representative SA-specific electrode (E1): tongue tip (TT) x-axis  
289 trajectories for /ba/, /da/, and /du/ (**f**) show distinct retraction magnitudes (marked by black arrow).  
290 Beta1 activity differs significantly across syllables during SA (**g**, black dots indicate time points  
291 where  $p < 0.05$ , one-way ANOVA) but not SI. SA beta1 activity is significantly higher than SI for all  
292 syllables (**h**, black dots indicate time points where  $p < 0.01$ , two-sided independent samples  $t$ -test).  
293 **i-k**, Representative SI-specific electrode (E2): tongue tip y-axis trajectories for /da/ and /du/  
294 highlight unique elevations (marked by black arrows) compared to /gu/ (**i**). Beta1 activity differs  
295 significantly during SI (**j**, black dots:  $p < 0.05$ , one-way ANOVA) but not SA. SI beta1 activity is  
296 significantly higher than SA for /da/ and /du/ (**k**, black dots:  $p < 0.01$ , two-sided independent  
297 samples  $t$ -test), while no significant difference is observed for /gu/. **l-n**, Representative SA-SI  
298 shared encoding electrode (E3, Cluster 1): lower incisor y-axis trajectories (**l**) show varying  
299 degrees of deliberate jaw elevation (indicated by black arrows) for /ba/, /da/, and /du/. Beta1 activity  
300 differs across syllables in both tasks (**m**, black dots:  $p < 0.05$ , one-way ANOVA) but shows no  
301 significant differences between modalities (**n**). Instead, a strong correlation is observed between  
302 tasks (Pearson's  $r$  and  $p$  values annotated). **o-r**, Representative Cluster 2 electrode (E4): modality-  
303 specific encoding electrode. SA-specific beta1 activity (**p**) encodes unique tongue tip retraction (**o**,  
304 marked by black arrows) for /da/ and /du/ compared to /bu/, while SI-specific beta1 activity (**r**)  
305 highlights distinct upper lip closure (**q**, marked by black arrows) for /bu/ compared to /da/ and /du/.  
306 Significant differences were marked (black dots:  $p < 0.05$ , one-way ANOVA). Time 0 ms indicates  
307 syllable onset (SA) or inferred onset (SI). Solid lines: mean; shaded areas: s.e.m. across syllable  
308 repetitions (color-coded). **s-t**, **w**, AKT encoding patterns derived from the temporal receptive field  
309 (TRF) model for SA-SI shared encoding electrodes during SA (**s**) and SI (**t**) tasks. Electrodes in (**s**)  
310 are organized by hierarchical clustering; the same order is maintained in (**t**). (**w**) Histogram of  
311 permutation testing demonstrates significant AKT encoding pattern similarity between SA-SI  
312 shared encoding electrodes during SA and SI tasks (red dashed line, *cosine similarity* = 0.28,  $p <$   
313 0.0001). **u-v**, **x**, AKT encoding patterns of SA-specific (**u**) and SI-specific (**v**) electrodes show no  
314 significant similarity ( $p = 0.0674$ , **x**) based on permutation testing (red dashed line, *cosine similarity*  
315 = 0.16,  $p = 0.0670$ ). The color intensity in (**s-v**) reflects the unique  $R^2$  to the full TRF models  
316 (scaled to 0-1 a.u.).

317

318 **Spatiotemporal hierarchy separates supramodal and modality-specific articulatory**  
319 **representations**

320 To investigate the spatial distribution patterns of SA-SI shared, SA-specific, and SI-specific  
321 encoding electrodes at the population level, we normalized all encoding electrodes to the MNI  
322 152 template and applied kernel density estimation (KDE) to construct probability density  
323 distributions for each encoding category (**Fig. 3a–c**). SA-specific encoding electrodes (**Fig. 3a,**  
324 **e**) and SI-specific encoding electrodes (**Fig. 3c, h**) were concentrated bilaterally around the  
325 central sulcus, predominantly covering the ventral primary sensorimotor cortex (areas 4, 3b, 1,  
326 and 2, according to a multi-modal parcellation atlas <sup>36</sup>). The peak point for SA-specific  
327 electrodes was located posteroinferior to that of SI-specific electrodes (**Fig. 3m**).

328 In contrast, supramodal (SA-SI shared encoding) electrodes (**Fig. 3b**) exhibited three  
329 distributed clusters (**Fig. 3b**): (1) middle premotor cortex (spanning areas 55b and dorsal 6v  
330 with peak density in dorsal 6v), (2) subcentral gyrus (encompassing areas 43 and ventral 6v  
331 centered on area 43), and (3) supramarginal-postcentral junction (covering areas OP4, PFop,  
332 and PF with focal concentration in OP4). These three peaks of supramodal electrodes were  
333 positioned anterosuperior, anteroinferior, and posteroinferior relative to those of SA-specific  
334 and SI-specific electrodes (**Fig. 3d, g**), exhibiting y-axis displacement (two-sided Kolmogorov-  
335 Smirnov test:  $p < 0.001$  vs SA-specific;  $p = 0.052$  vs SI-specific) (**Fig. 3d–i**).

336



341 electrodes (purple, **b**), and SI-specific encoding electrodes (blue, **c**) using kernel density estimation  
342 (KDE). Darker colors indicate higher probability density, with a cut-off at 0.5 cumulative density  
343 (a.u.). Yellow stars mark the peak density points, and black boxes indicate the regions shown in  
344 panels **d**, **g**, and **m**. **d**, KDE contour plot comparing SA-specific (red) and SA-SI shared (purple)  
345 encoding electrodes during the SA task. KDE curves along the Y- and Z-axes of MNI space are  
346 shown, with  $p$ -values from two-sided Kolmogorov-Smirnov tests. The gray solid lines mark key  
347 brain sulci: prCS = precentral sulcus, CS = central sulcus, poCS = postcentral sulcus, SF = Sylvian  
348 fissure, SFS = superior frontal sulcus, IFS = inferior frontal sulcus, IPS = intraparietal sulcus. **e-f**,  
349 Spatial scatter plots of SA-specific encoding electrodes (**e**, red) and SA-SI shared encoding  
350 electrodes (**f**, purple) during the SA task. The size of each point reflects the  $R^2$  value in the SA-  
351 AKT model. **(g)** KDE contour plot comparing SI-specific (blue) and SA-SI shared (purple) encoding  
352 electrodes during the SI task. KDE curves are displayed along the Y- and Z-axes with  $p$ -values  
353 from two-sided Kolmogorov-Smirnov tests. **h-i**, Spatial scatter plots of SI-specific encoding  
354 electrodes (**h**, blue) and SA-SI shared encoding electrodes (**i**, purple) during the SI task. Point size  
355 reflects the  $R^2$  value in the SI-AKT model. **j**, Violin plots illustrating the onset (top) and peak  
356 (bottom) times of SA-SI shared encoding electrodes during the SA task (purple) and SI task (light  
357 purple), showing no significant differences based on two-sided Mann-Whitney  $U$  tests. **k**, Violin  
358 plots showing the onset (top) and peak (bottom) times of SA-specific (red) and SA-SI shared  
359 (purple) encoding electrodes during the SA task, with later timings for SA-specific electrodes based  
360 on two-sided Mann-Whitney  $U$  tests. **l**, Violin plots showing the onset (top) and peak (bottom) times  
361 of SI-specific (blue) and SA-SI shared (purple) encoding electrodes during the SI task, with no  
362 significant differences observed based on two-sided Mann-Whitney  $U$  tests. **m**, KDE contour plot  
363 comparing the spatial distribution of SA-specific (red) and SI-specific (blue) encoding electrodes.  
364 KDE curves are displayed along the Y- and Z-axes with  $p$ -values from two-sided Kolmogorov-  
365 Smirnov tests. **n**, Violin plots comparing the onset (left) and peak times (right) of SA-specific (red)  
366 and SI-specific (blue) encoding electrodes during their respective tasks, with later timings observed  
367 for SA-specific electrodes based on two-sided Mann-Whitney  $U$  tests.

368

369 Next, we analyzed the temporal dynamics of the different electrode groups, specifically  
370 comparing across electrode groups for the onset and peak times in SA and SI tasks. Our analysis  
371 revealed distinct temporal activation patterns across neural populations (**Fig. 3j-n**).  
372 Supramodal (SA-SI shared encoding) electrodes showed remarkable temporal stability, with  
373 no significant difference in median onset time (median [IQR], SA: 480 ms [328–710 ms], SI:  
374 480 ms [340–570 ms],  $p = 0.187$ , two-sided Mann-Whitney  $U$  test) or peak time (SA: 810 ms  
375 [590–980 ms], SI: 720 ms [640–885 ms],  $p = 0.158$ ) between modalities (**Fig. 3j**). During SA  
376 task, these supramodal electrodes activated significantly earlier than SA-specific electrodes  
377 (onset: 580 ms [235–820 ms], 100 ms difference,  $p = 0.049$ ; peak: 940 ms [680–1120 ms], 130  
378 ms difference,  $p < 0.001$ ), forming a temporally cascaded processing stream (**Fig. 3k**). During  
379 the SI task, SI-specific electrodes showed synchronous activation (onset: 460 ms [235–590 ms],  
380  $p = 0.379$ ; peak: 720 ms [600–880 ms],  $p = 0.640$ ; **Fig. 3l**) with supramodal electrodes. In  
381 contrast, SA-specific electrodes exhibited consistently delayed activation compared to SI-  
382 specific electrodes across all temporal measures (all  $p < 0.001$ ), revealing fundamentally  
383 different processing architectures between production modalities (**Fig. 3n**).

384 These findings suggest a hierarchical processing model where supramodal populations in  
385 middle premotor cortex, subcentral gyrus, and supramarginal-postcentral junction represent a  
386 common upstream for abstract articulatory planning. During overt speech articulation, this  
387 supramodal activity cascades to SA-specific populations in the ventral primary sensorimotor  
388 cortex for motor execution. During mental imagery, concurrent SI-specific activation across  
389 the same regions preserves internal articulatory output, enabling purely cognitive speech  
390 simulation while generating quasi-perceptual experience <sup>6</sup>.

391

### 392 **Supramodal integration versus modality-specific somatotopic organization of** 393 **articulatory gestures**

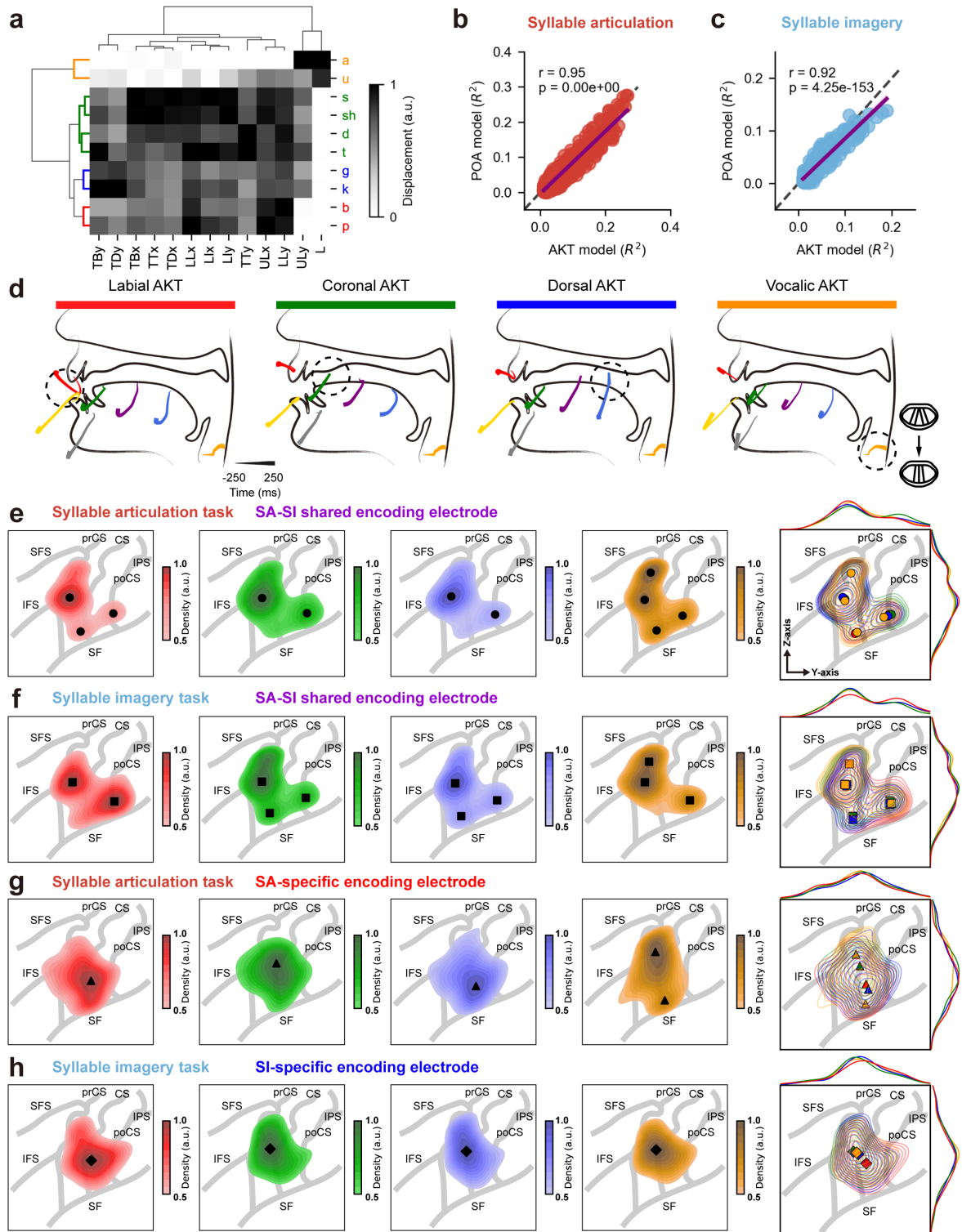
394 Previous ECoG studies have revealed that neural dynamics at the single-electrode level encode  
395 a diverse repertoire of AKTs during continuous speech production <sup>28</sup>, each reflecting  
396 coordinated movements of articulators toward four primary vocal tract shapes corresponding

397 to distinct places of articulation (POAs): coronal, labial, dorsal constrictions, and vocalic  
398 control. Based on this, we hypothesized that beta1 activity in SA and SI modalities exhibits  
399 similar encoding modes.

400 To test this hypothesis, we first examined whether AKTs associated with different  
401 phonemes could be organized according to the similarity of their articulatory gestures. We  
402 computed the displacement matrix of AKTs during articulation and performed unsupervised  
403 hierarchical clustering across phonemes. This analysis revealed a clear grouping of phonemes  
404 into four established POA categories, with subtle intra-cluster variations (**Fig. 4a**). We then  
405 aggregated AKTs within each POA-defined cluster and found that core articulatory features  
406 were preserved, reinforcing the consistency of gesture patterns within each category (**Fig. 4d**).

407 To further evaluate the encoding of these POA features in beta1 activity, we trained a TRF  
408 model incorporating POA cluster-based features, marking the presence or absence of each  
409 phoneme's POA feature at its onset using binary indicators (1/0). We found a strong positive  
410 correlation between the total variance explained ( $R^2$ ) of the 13-dimensional AKT-based TRF  
411 model and the 4-dimensional POA-based TRF model in both SA and SI modalities (SA:  
412 Pearson's  $r = 0.95$ ,  $p < 0.0001$ ; SI:  $r = 0.92$ ,  $p < 0.0001$ ) (**Fig. 4b-c**), indicating a high degree  
413 of consistency in the neural encoding of AKT and POA representations across both modalities.  
414 This also suggests that the frontoparietal regions exhibit a relatively rigid encoding of AKT  
415 combinations, which closely aligns with the four distinct articulatory gestures.

416



417

418 **Fig. 4 | Somatotopic arrangement of four distinct vocal tract gestures in SA-specific, SI-**  
 419 **specific, and SA-SI shared encoding electrodes.**

420 **a**, Hierarchical clustering of phonemes based on AKT displacement follows the place of articulation  
 421 (POA) categories (red: labial, green: coronal, blue: dorsal, orange: vocalic). **b-c**, Scatter plots show  
 422 strong correlations between the total variance explained ( $R^2$ ) of the 13-dimensional AKT and 4-

423 dimensional POA-based TRF models during both syllable articulation (**b**) and imagery (**c**) tasks  
424 (SA: Pearson's  $r = 0.95$ ,  $p < 0.0001$ ; SI: Pearson's  $r = 0.92$ ,  $p < 0.0001$ ), with purple fitted curves  
425 closely aligned to the  $y = x$  diagonal (black dashed line). **d**, Averaged movement trajectories of  
426 seven articulators grouped by four categories of POA during phoneme articulation across all nine  
427 subjects. The time window spans from -250 ms to 250 ms relative to phoneme onset. The time  
428 course of trajectories is represented by thin-to-thick lines. Dashed circles highlight the primary  
429 vocal tract constriction for each POA category: labial AKT (red) involves lip closure, coronal AKT  
430 (green) involves tongue tip elevation contacting the alveolar ridge, dorsal AKT (blue) involves  
431 tongue dorsum elevation against the soft palate, and vocalic AKT (orange) involves the constriction  
432 and approximation of the vocal folds for voicing. **e-h**, Weighted KDE probability density  
433 distributions for the four POA categories during the syllable articulation task (**e, g**) and the syllable  
434 imagery task (**f, h**), visualized for SA-SI shared (**e, f**), SA-specific (**g**), and SI-specific (**h**) encoding  
435 electrodes. In the first four columns of panels **e-h**, the distributions were weighted by the  $\Delta R^2$   
436 values (normalized to a scale of 0–1) of the POAs, with a cut-off at 0.5 cumulative density (a.u.).  
437 Each color represents a distinct POA: red for labial, green for coronal, blue for dorsal, and yellow  
438 for vocalic. The intensity of the color corresponds to the probability density, with darker shades  
439 indicating higher density. Circles (**e**), squares (**f**), triangles (**g**), or diamonds (**h**) indicate the  
440 positions of the probability density peaks. The last column of panels **e-h** displays joint weighted  
441 KDE contour plots, with colored symbols marking the peak locations for each POA. Gray solid lines  
442 indicate major brain sulci: prCS = precentral sulcus, CS = central sulcus, poCS = postcentral sulcus,  
443 SF = Sylvian fissure, SFS = superior frontal sulcus, IFS = inferior frontal sulcus, and IPS =  
444 intraparietal sulcus. Refer to **Extended Data Figs. 6 and 7** for original scatter plots, all peak  
445 locations, and pairwise permutation test results of POA encoding across SA-specific, SI-specific,  
446 and SA-SI shared encoding electrodes during SA and SI modalities.

447

448

449 Subsequently, we aimed to determine the somatotopic arrangements across POAs for SA-  
450 SI shared, SA-specific, and SI-specific encoding electrodes, respectively. To this end, we  
451 constructed weighted functional maps for each POA in the y-z plane of MNI space, using each  
452 POA's unique contribution ( $\Delta R^2$ ) as the weighting factor in the KDE to visualize the spatial  
453 encoding patterns for the three categories of electrodes.

454 Our spatial analysis revealed fundamentally distinct organizational principles between  
455 supramodal and modality-specific articulatory representations. For SA-SI shared encoding  
456 electrodes, we observed overlapping functional peaks for multiple POAs in three key regions:  
457 the middle premotor cortex, subcentral gyrus, and POCG-SMG junction (**Fig. 4e-f**, consistent  
458 with peak locations identified in **Fig. 3b**). This conserved spatial patterning across both speech  
459 modalities (*cosine similarity* = 0.56,  $p < 0.0001$ , **Extended Data Fig. 5**) suggests that these  
460 supramodal regions implement integrative planning and coordination of multiple articulatory  
461 gestures<sup>10,37-39</sup>. In addition, we identified a functional peak for vocalic constriction at the MFG-  
462 PRCG junction under both SA and SI modalities, suggesting a preferential encoding for vocalic  
463 control planning in this region.

464 In contrast, modality-specific electrodes exhibited distinct somatotopic arrangements. SA-  
465 specific electrodes exhibited a well-defined ventral-dorsal somatotopic gradient in the vSMC,  
466 with dual vocalic representations at the ventral and dorsal extremes. Between these vocalic  
467 peaks, we observed sequential representations of coronal, labial, and dorsal POAs, with  
468 significant separation between all adjacent peaks (**Fig. 4g**; pairwise permutation test results  
469 shown in **Extended Data Fig. 6**). Furthermore, SI-specific encoding electrodes displayed a  
470 less distinct somatotopic arrangement along the central sulcus, oriented from superior-anterior  
471 to inferior-posterior, with coronal, vocalic, dorsal, and labial peaks positioned sequentially (**Fig.**  
472 **4h**, **Extended Data Fig. 6**). However, only the separation between labial and coronal peaks  
473 reached statistical significance ( $p = 0.0336$  along the MNI y-axis; pairwise permutation tests,  
474 Benjamini-Hochberg corrected, **Extended Data Fig. 6**).

475 Our findings reveal a hierarchical somato-cognitive architecture for speech imagery and  
476 articulation. Within distributed frontoparietal nodes, spatially condensed supramodal  
477 assemblies support compressed articulatory gesture planning, enabling computational

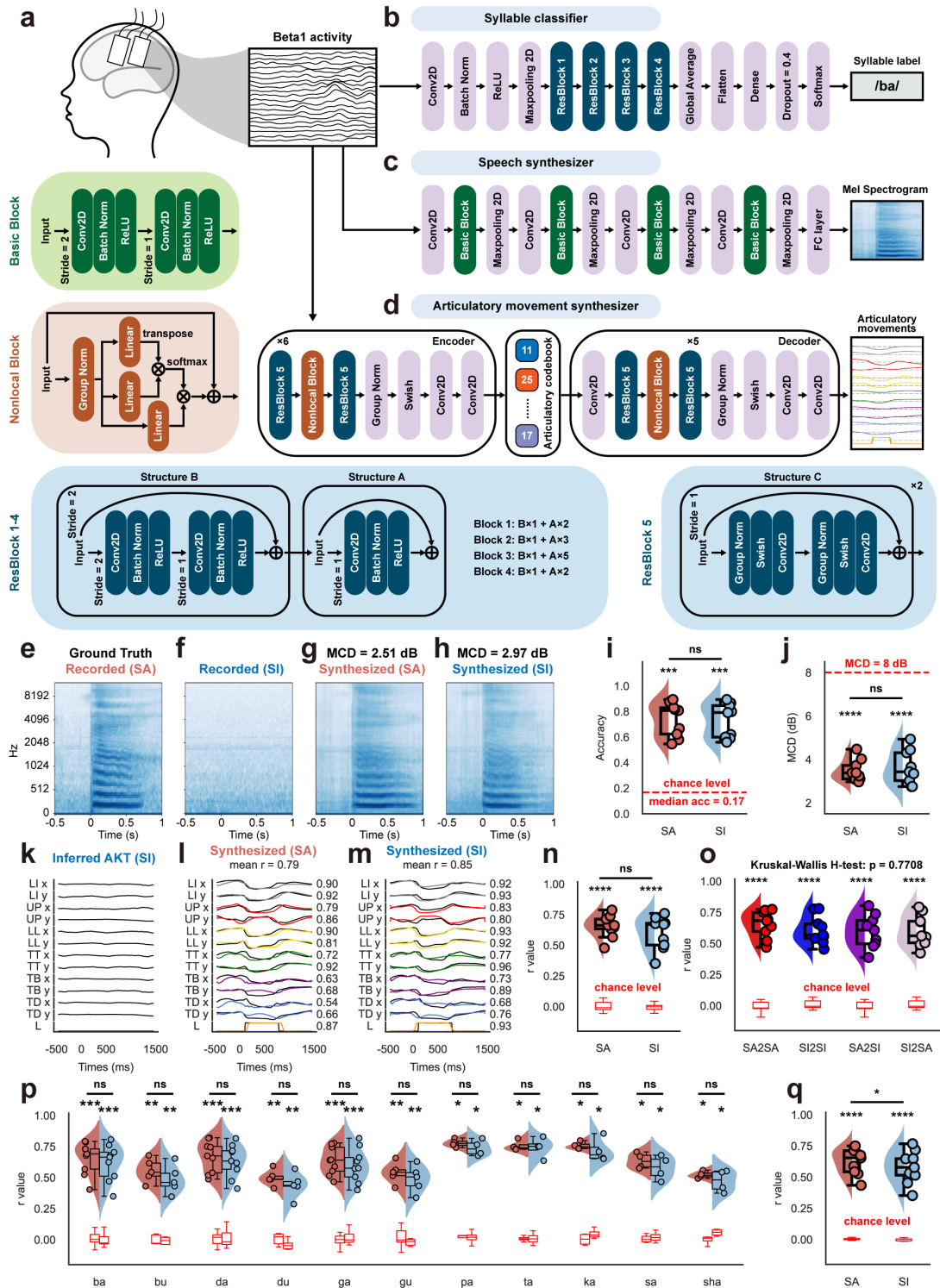
478 efficiency. These abstract representations dynamically interface with modality-specific  
479 somatotopic implementations in primary sensorimotor cortex: the system maintains clear  
480 somatotopic organization during overt articulation, whereas covert imagery engages a  
481 reconfigured functional topography that facilitates internal simulation.

482

483 **Frontoparietal articulatory codes enable effective syllable classification, speech synthesis,**  
484 **and kinematic decoding across executed and imagined speech**

485 Having characterized the encoding patterns of AKTs and POAs at the single-electrode level for  
486 both modalities, we next asked whether population-level articulatory codes in frontoparietal  
487 cortex could support speech decoding, particularly for inner speech in the SI condition. To  
488 address this, we developed a triple-stream neural decoding framework (**Fig. 5a**) that maps beta1  
489 activity from SA- and SI-responsive electrodes to three parallel outputs: syllable classification  
490 (syllable classifier; **Fig. 5b**), acoustic waveform reconstruction (speech synthesizer; **Fig. 5c**),  
491 and AKTs (articulatory movement synthesizer; **Fig. 5d**). Decoder performance was evaluated  
492 across all nine subjects using syllable classification accuracy, Mel cepstral distortion (MCD)  
493 for synthesized audio, and Pearson's correlation coefficient ( $r$ ) for predicted articulatory  
494 trajectories.

495



496 **Fig. 5 | Deep learning-based synthesis of text, speech, and articulatory kinematics from**  
 497 **frontoparietal beta1 activity during SA and SI tasks.**

498 **a**, Schematic diagram of the deep learning model pipeline for synthesizing text, speech and  
 499 AKTs from frontoparietal beta1 activity. **b-d**, Model architectures for the syllable classifier (**b**), speech  
 500 audio synthesizer (**c**), and articulatory movement synthesizer (**d**). **e**, Recorded Mel-spectrogram  
 501 of a /ba/ during the SA task (ground truth). **f**, Recorded Mel-spectrogram during syllable imagery  
 502 of a /ba/ (blank control). **g-h**, Mel-spectrogram synthesized from beta1 activity recorded during SA

503 (g), and SI (h) modalities, with Mel-Cepstral Distortion (MCD) values of 2.51 dB and 2.97 dB,  
504 respectively. i, Violin plots of syllable classifier accuracy for electrodes responsive to SA and SI  
505 tasks (subject  $N = 9$ ). The red dashed line indicates the median chance level; \*\*\* $p < 0.001$  for  
506 Mann–Whitney  $U$  test against chance level; ns: non-significant, Wilcoxon signed-rank test for SA  
507 vs. SI. j, Violin plot of synthesized speech MCD values. The red dashed line represents the 8 dB  
508 threshold considered acceptable for voice recognition; \*\*\*\* $p < 0.0001$  for Mann–Whitney  $U$  test  
509 against the threshold; ns: non-significant, Wilcoxon signed-rank test for SA vs. SI. k, AKTs inferred  
510 from audio recorded during SI modality using AAI algorithm (blank control). l-m, Synthesized AKTs  
511 from beta1 activity following Leave-One-Syllable-Out strategy in SA (l) and SI (m) tasks, with mean  
512 Pearson's  $r$  values across all AKTs of 0.79 and 0.85, respectively. Black lines indicate ground truth  
513 AKTs inferred from audio recorded during SA modality. n, Violin plot comparing the performance  
514 of synthesized AKTs, represented by mean Pearson's  $r$  values of all AKTs with ground truth, for SA  
515 and SI-responsive electrodes (subject  $N = 9$ ). The red box represents chance level; \*\*\*\* $p < 0.0001$   
516 for Mann–Whitney  $U$  test against chance level; ns: non-significant, Wilcoxon signed-rank test for  
517 SA vs. SI. o, Violin plot of cross-modal transfer learning performance for SA-SI shared electrodes  
518 (subject  $N = 9$ ), with the red box marking chance level. X-axis labels are formatted as "training set  
519 to test set" (e.g., SA2SI indicates a model trained on SA and tested on SI); \*\*\*\* $p < 0.0001$  for  
520 Mann–Whitney  $U$  test against chance level. Kruskal-Wallis H-test across models:  $p = 0.7708$ . p,  
521 Leave-One-Syllable-Out performance of all syllables for SA and SI responsive electrodes. The x-  
522 axis represents the syllables left out. The red box represents chance level; \*, \*\*, \*\*\*, \*\*\*\* denote  
523 Mann–Whitney  $U$  test significance at  $p < 0.05$ , 0.01, 0.001, and 0.0001, respectively, against  
524 chance level; ns: non-significant, Wilcoxon signed-rank test for SA vs. SI. q, Summary violin plot  
525 for Leave-One-Syllable-Out data comparing SA and SI modalities (subject  $N = 9$ ); \*\*\*\* $p < 0.0001$   
526 for Mann–Whitney  $U$  test against chance level, \* $p < 0.05$  (Wilcoxon signed-rank test for SA vs. SI).  
527 All data in the violin and box plots represent median and interquartile ranges across subjects.

528

529 The syllable classifier achieved high decoding accuracy in both SA and SI tasks, with  
530 comparable performance (SI: median accuracy = 0.79, IQR = 0.60–0.85; SA: 0.81, IQR = 0.62–  
531 0.83;  $p = 0.3594$ , Wilcoxon signed-rank test). Crucially, decoding in the SI modality was  
532 significantly above chance (median chance = 0.17,  $p = 0.0003$ , Mann-Whitney  $U$  test),  
533 indicating that internal articulatory representations are sufficiently stable to support categorical  
534 decoding (**Fig. 5i**).

535 Given that no audio is generated in the SI modality (as shown in **Fig. 5f** for covertly reading  
536 /ba/), we used the audio recorded in the SA modality for the same trial within the same block  
537 as the shared ground truth (see **Fig. 5e** for overtly reading /ba/) for both the SA and SI  
538 modalities to train the speech synthesizer. The Mel cepstral distortion (MCD) between the  
539 synthesized and real speech spectrograms was significantly below the 8 dB threshold in both  
540 the SA (median MCD = 3.38 dB, IQR: 3.11–3.72 dB) and SI (median MCD = 3.43 dB, IQR:  
541 3.03–4.31 dB) modalities (both  $p = 0.0039$ ), which is considered acceptable for voice  
542 recognition systems<sup>30</sup>. SI decoding achieved speech synthesis quality comparable to SA ( $p =$   
543 0.3008), supporting the viability of mental imagery for high-fidelity output (**Fig. 5j**; see **Fig.**  
544 **5g-h** for synthesized /ba/ Mel spectrograms for both modalities).

545 Similarly, using inferred articulatory trajectories from SA trials as the ground truth, we  
546 decoded AKTs from beta1 activity using all SI-responsive electrodes (**Fig. 5k–m**). Synthesized  
547 trajectories in the SI modality showed strong correlations with ground truth (median  $r = 0.67$ ,  
548 IQR = 0.50–0.68), comparable to those observed in SA (0.66, IQR = 0.63–0.71;  $p = 0.0547$ ),  
549 and both well above chance ( $p = 0.0004$ ) (**Fig. 5n**). Similar decoding performance was  
550 observed when using only SI-specific electrodes (**Extended Data Fig. 8**), further supporting  
551 that internal articulatory dynamics are robustly encoded in SI-related beta1 activity.

552 As a result, deep neural network models trained and tested within each task demonstrated  
553 equivalent decoding efficacy for SA and SI tasks. Next, we investigated whether cross-modal  
554 articulatory movement decoding could be achieved using SA-SI shared electrodes. We trained  
555 four articulatory movement synthesizers based on these electrodes: SA2SA, SA2SI, SI2SI, and  
556 SI2SA (naming format indicates "training set to testing set," e.g., SA2SI represents training  
557 with SA modality beta1 activity and testing with SI modality beta1 activity, evaluating the

588 correlation between synthesized movement and ground truth). Overall, we found no significant  
589 difference in the median correlation coefficients (SA2SA: 0.68 [0.60–0.75], SA2SI: 0.61  
590 [0.50–0.69], SI2SI: 0.57 [0.54–0.66], SI2SA: 0.56 [0.51–0.68],  $p = 0.7708$ , Kruskal-Wallis  $H$ -  
591 test), with all models significantly above chance level (all  $p = 0.0004$ , Mann-Whitney  $U$  test)  
592 (**Fig. 5o**). These results further support that SA-SI shared electrodes encode supramodal  
593 articulatory representations, offering a principled basis and implantation target for cross-modal  
594 speech BCIs. Such systems may be especially beneficial for individuals with progressively  
595 declining speech function, such as those diagnosed with ALS.

596 We next aimed to validate whether our articulatory movement synthesizer had effectively  
597 learned the stable relationship between beta1 activity and AKTs. If this were the case, the  
598 synthesizer should exhibit extrapolative generalization, meaning it should be capable of  
599 predicting the movement trajectories of syllables not included in the training set. To assess this,  
600 we employed a Leave-One-Syllable-Out strategy, iteratively excluding one syllable, training  
601 the synthesizer on the remaining syllables, and testing the excluded syllable to compute the  
602 correlation ( $r$ ) between predicted AKTs and ground truth. We found that the Leave-One-  
603 Syllable-Out articulatory movement synthesizer still performed well for all syllables in both  
604 SA and SI responsive electrodes, with the median correlation coefficient for all syllables  
605 significantly exceeding chance levels ( $p < 0.05$  for all syllables, Mann-Whitney  $U$  test) (**Fig.**  
606 **5p**; see **Fig. 5l-m** for synthesized articulatory movements of /ba/ in SA and SI modalities using  
607 Leave-One-Syllable-Out strategy). In summary, the Leave-One-Syllable-Out model  
608 demonstrated strong performance in both the SA (median  $r = 0.65$ , IQR: 0.54–0.71) and SI  
609 (median  $r = 0.58$ , IQR: 0.52–0.65) modalities (all  $p = 0.0004$  vs chance level, **Fig. 5q**).  
610 Collectively, our articulatory movement synthesizer demonstrated robust extrapolative  
611 generalization across both modalities, indicating that it successfully captured a  
612 neurocomputationally interpretable mapping between cortical signatures and intended  
613 articulatory dynamics.

584

585

## 586 **Discussion**

587 Our high-density ECoG recordings reveal a hierarchical cortical architecture underpinning  
588 shared and modality-specific mechanisms for articulatory control in speech imagery and  
589 articulation. Three key advances emerge: First, frontoparietal AKT encoding exhibits a spectral  
590 dissociation, with high-gamma dominance and secondary beta1 contributions in SA,  
591 contrasting with a primary beta1 signature in SI. Second, beta1-band analysis identifies  
592 spatiotemporally stable supramodal populations in the middle premotor cortex, subcentral  
593 region, and postcentral–supramarginal areas that support integrated multi-articulatory planning.  
594 These populations hierarchically interface with somatotopically interleaved primary  
595 sensorimotor populations to segregate internal simulation during speech imagery from external  
596 execution during articulation. Third, neural decoding based on these neural dynamics enables  
597 robust within-modal classification, speech synthesis, and kinematic prediction, while also  
598 demonstrating cross-modal generalization and extrapolation to untrained syllables in both  
599 modalities.

600

### 601 **Distinct spectral signatures across speech modalities**

602 Our findings reveal a fundamental dissociation in spectral encoding strategies between speech  
603 imagery and articulation. High-gamma activity emerges as the dominant carrier of AKT  
604 information during overt speech articulation, particularly in frontoparietal regions, consistent  
605 with its established role in encoding fine-grained articulatory features, including AKT  
606 sequences, place of articulation, and phonetic details<sup>19-21,25,27,28,40,41</sup>. Notably, while residual  
607 high-gamma activity persists during speech imagery<sup>19-21,40</sup>, its functional engagement is  
608 markedly attenuated—both in spatial extent (reduced electrode recruitment) and decoding  
609 performance ( $R^2$  ranked fourth among frequency bands).

610 Conversely, beta1 activity exhibits maximal responsiveness and reliably encodes AKTs  
611 across both speech modalities. This band not only achieves secondary decoding performance  
612 during articulation (surpassed only by high gamma) but also emerges as the principal  
613 information carrier during imagery throughout frontoparietal regions. This dual-domain

614 involvement aligns with converging evidence from human and primate studies that implicate  
615 beta1 activity in motor planning and execution, action simulation and imagery, and postural  
616 maintenance in humans and non-human primates<sup>40,42-47</sup>.

617

### 618 **An expanded somato-cognitive architecture for speech imagery and articulation**

619 Recent work has reconceptualized the sensorimotor cortex as a somato-cognitive action  
620 network (SCAN), comprising interdigitated effector-specific zones and integrative inter-  
621 effector regions<sup>39,48</sup>. By focusing on beta1 activity, our findings extend this ‘new homunculus’  
622 framework by revealing functional distinctions between supramodal and modality-specific  
623 populations involved in speech planning, internal simulation and external execution.

624 In addition to previously described inter-effector middle premotor and subcentral hubs<sup>39,48</sup>,  
625 we identify a key parietal node at the POCG–SMG junction that, together with frontal  
626 counterparts, forms an extended fronto-parietal network for articulatory planning. Our data  
627 reveals that supramodal populations within these hubs exhibit three key features: robust cross-  
628 modal stability across spatiotemporal profiles and AKT encoding patterns; spatially  
629 overlapping encoding of multiple articulatory gestures, supporting integrative control; and  
630 consistent temporal precedence over downstream somatotopically organized, modality-specific  
631 populations—by ~100 ms at onset and ~130 ms at peak activity. These functional signatures  
632 are aligned with direct cortical stimulation studies, which show that disruption of these peri-  
633 Sylvian regions reliably induces complete speech arrest, implicating their role in coordinated  
634 articulatory planning rather than isolated motor output<sup>38,49</sup>.

635 Further converging evidence comes from fMRI studies showing strong connectivity  
636 between these hubs and the cingulo-opercular network (CON), a domain-general system  
637 implicated in goal-directed action and performance monitoring<sup>39,50,51</sup>, as well as from  
638 anatomical studies showing their tight coupling via the superior longitudinal fasciculus-III /  
639 arcuate fasciculus complex—critical tracts within the dorsal phonological stream<sup>9,36,49,52-54</sup>.  
640 Stimulation of these fiber tracts likewise disrupts speech, resulting in phonemic paraphasia or  
641 complete speech arrest<sup>52,55-57</sup>. Utah array recordings have further shown that limited coverage

642 within these hubs suffices for decoding natural speech involving complex multi-articulator  
643 coordination<sup>26,58,59</sup>. Together, these findings support the interpretation that supramodal  
644 populations encode cognitive-level articulatory plans generalizable across modalities. This  
645 provides direct neurophysiological support for the motor simulation hypothesis over the  
646 abstraction view, suggesting that speech imagery entails precise, temporally structured  
647 articulatory planning that is initiated but not executed<sup>6,11,13</sup>.

648 In contrast, modality-specific populations show distinct anatomical and functional profiles.  
649 Predominantly localized to the primary sensorimotor cortex along the central sulcus, SA- and  
650 SI-specific populations are interleaved in a mosaic-like distribution that, while spatially  
651 overlapping, diverges in temporal dynamics and cross-modal encoding patterns.

652 Consistent with prior findings, SA-specific encoding populations exhibited a characteristic  
653 somatotopic organization, in sharp contrast to the integrative arrangement of supramodal  
654 regions. Vocalic AKTs evoked two spatial peaks at the dorsal and ventral extremes of vSMC,  
655 with coronal, labial, and dorsal peaks arranged sequentially along the superior–inferior axis.  
656 These vocalic peaks align closely with laryngeal motor areas identified in previous ECoG and  
657 fMRI work<sup>27,60</sup>. Cortical stimulation at these vocalic sites induces involuntary vocalizations,  
658 pitch change, laryngeal constriction, or paralysis, providing causal evidence for their functional  
659 specificity<sup>29,48,61-63</sup>. Temporally, the activation of supramodal populations consistently  
660 preceded that of SA-specific populations, echoing the ~145 ms gap between phonetic encoding  
661 and articulatory output observed in large-scale neuroimaging meta-analyses<sup>64-66</sup>.

662 In contrast, SI-specific electrodes exhibited a remapped functional organization with  
663 attenuated somatotopic segregation. This reduced somatotopic differentiation likely stems from  
664 their sparsely distributed and spatially dispersed organization. This may also explain the limited  
665 detection of SI-specific activity in primary sensorimotor cortex in prior fMRI studies<sup>67,68</sup>. The  
666 dominant motor signal during articulation likely masks more subtle SI-related activations in  
667 cross-modal contrasts.

668 Remarkably, our high-density ECoG recordings identified distinct populations that  
669 exclusively encode internal AKTs during speech imagery (e.g., **Fig. 2i–k**), providing direct

670 neurophysiological evidence for cognitive-level articulatory outputs. These internal  
671 simulations paralleled the temporal profiles of supramodal populations and may contribute to  
672 the quasi-perceptual experience of speech<sup>17,18</sup>. At the same time, they may suppress overt  
673 execution via local inhibitory interactions with neighboring SA-specific circuits, enabling  
674 mental speech simulation without triggering actual articulation<sup>7,10,17,18,69-71</sup>.

675 Together, these findings extend the SCAN framework to a broader fronto-parietal  
676 architecture linking overt articulation and covert imagery. Within this structure, higher-level,  
677 supramodal articulatory planning regions interface with modality-specific execution systems,  
678 forming a hierarchical somato-cognitive map in the sensorimotor cortex. This hierarchy is  
679 mirrored by anatomical gradients in myelin content, receptor density, and cytoarchitectonic  
680 differentiation, underscoring the specialized contributions of these regions to speech control  
681 across internal and external domains<sup>36,72-74</sup>.

682

### 683 **Implications for neurally interpretable within- and cross-modal speech BCIs**

684 Our encoding results establish the basis for a biologically grounded framework for decoding  
685 both spoken and imagined speech. Despite the constraints of a limited ECoG dataset, we  
686 developed a triple-stream architecture that accurately classified syllables, synthesized acoustic  
687 output, and predicted AKTs by leveraging frontoparietal beta1 activity—an underutilized  
688 neural signal in prior BCI efforts. Within each modality, this framework achieved performance  
689 levels comparable to those reported in intraoperative and early-stage implanted BCI studies  
690<sup>20,21,23,58,59,75</sup>, with considerable potential for further improvement<sup>25,26,76</sup>, particularly in  
691 decoding speech imagery.

692 A key advance underlying this framework is the discovery of supramodal articulatory  
693 representations within the frontoparietal cortex. These representations exhibited robust cross-  
694 modal encoding stability, enabling generalization between modalities with prediction  
695 consistency approaching within-modality levels. This degree of cross-modal transferability  
696 suggests a shared speech planning architecture, capable of integrating both overt articulation  
697 and covert speech imagery. Such systems may be particularly valuable for individuals with

698 progressive speech impairment, including those affected by ALS. Critically, the localization of  
699 these supramodal representations to the middle premotor cortex, subcentral gyrus, and POCG–  
700 SMG junction provides anatomically and functionally grounded targets for future cross-modal  
701 BCI implantation. More broadly, by engaging higher-order speech planning nodes (e.g., the  
702 SMG node) and bypassing downstream corticobulbar pathways, this approach may extend  
703 beyond brainstem or lower tract dysfunction to encompass cortical impairments such as  
704 Broca’s aphasia, thereby broadening the clinical applicability of speech BCIs.

705 Furthermore, the model generalized effectively to untrained syllables, demonstrating the  
706 utility of AKTs as a stable and interpretable intermediary between neural activity and phonemic  
707 categories. This feature not only enhances decoding transparency but also enables scalability  
708 beyond limited-vocabulary settings, a critical step toward open-set, real-world speech  
709 neuroprosthetics <sup>25,26</sup>.

710 Together, these results highlight a principled path toward within-modal and cross-modal  
711 speech BCIs across speech imagery and articulation. The integration of supramodal and  
712 modality-specific articulatory representations paves the way for high-performance, low-  
713 training neural prosthetics, expanding the clinical reach of speech restoration technologies for  
714 patients across a broad spectrum of speech output disorders, particularly those with locked-in  
715 syndrome, Broca’s aphasia, or progressive conditions such as ALS <sup>20-22,58</sup>.

716

## 717 **Methods**

### 718 **Participants**

719 A total of 9 subjects from Huashan Hospital participated in this study. All participants (median  
720 age [IQR]: 35 [32, 51] years; 4 males, 5 females; 9 left hemisphere) were eloquent brain tumor  
721 patients undergoing awake language mapping as part of their surgery. During the intraoperative  
722 language mapping, high-density electrode grids were temporarily placed onto the frontoparietal  
723 cortex to record local field potentials from the cortex, and the participants were instructed to  
724 perform the experimental tasks.

725 Subjects were asked to participate in the research study only if they were undergoing  
726 awake surgery with direct cortical stimulation as part of the regular clinical routine, meaning  
727 that this was deemed necessary for the safe resection of their tumor. Each participant was  
728 consented prior to the surgery, at which time it was explained in a transparent manner (as  
729 detailed in the IRB-approved written protocol/consent document) that the research task was for  
730 scientific purposes and would not directly impact their care. It was clearly articulated to each  
731 subject that participation in the research task was completely voluntary. The experimental  
732 protocol was approved by the Huashan Hospital Institutional Review Board of Fudan  
733 University (HIRB, KY2022-504). All participants gave their written informed consent prior to  
734 testing.

735

### 736 **Experiment paradigm**

737 Two distinct paradigms were employed in this study. Five subjects (S1–S5) completed  
738 Paradigm 1 (**Fig. 1a**), while four subjects (S6–S9) completed Paradigm 2 (**Fig. 1b**). Both  
739 paradigms were presented using E-Prime 2.0 (Psychology Software Tools, Inc.) and were  
740 temporally synchronized with ECoG and audio recordings.

741 Paradigm 1 (**Fig. 1a**) consisted of two parts. In the first part, subjects were presented with  
742 the instruction "Please Press the key" and performed an ipsilateral key press task, repeated 40  
743 times, serving as an ipsilateral hand motor baseline control. In the second part, a visual stimulus  
744 displaying a syllable (e.g., "ba") was presented on the screen for 2 seconds, followed by the  
745 instruction "Prepare to Read Aloud" for another 2 seconds. Subjects were then required to  
746 overtly produce the syllable five times (SA task) upon the appearance of a black cross symbol  
747 (go-cue) while simultaneously marking the onset of each production with a button press (with  
748 intervals of at least 1.2 s). Subsequently, the instruction "Prepare for imagery" was displayed  
749 for 2 seconds, after which subjects covertly imagined producing the syllable five times (SI task)  
750 upon the appearance of the black cross symbol (go-cue), again marking each imagined  
751 production onset with a button press (intervals > 1.2 s). The ipsilateral key press served as the  
752 reference time point for aligning neural, audio, and feature matrix data.

753 Paradigm 2 (**Fig. 1b**) introduced an auditory cueing component. Subjects first listened to  
754 a syllable auditory stimulus (e.g., "ga") repeated five times. The syllable auditory stimuli were  
755 extracted from recordings of the subjects' own voices, obtained during preoperative training.  
756 This was followed by the instruction "Prepare to Read Aloud", accompanied by a visual display  
757 of the syllable (2 s). Next, a gray cross symbol and a black cross symbol alternated every 1 s,  
758 and subjects were instructed to overtly produce the syllable each time the black cross appeared,  
759 repeating this process five times (SA task). The syllable imagery task followed the same  
760 procedure: after the instruction "Prepare for imagery" and the visual presentation of the syllable  
761 (2 s), subjects covertly imagined producing the syllable in sync with the black cross symbol,  
762 repeating the task five times (SI task). The appearance of the black cross (go-cue) served as the  
763 reference time point for subsequent neural, audio, and feature matrix segmentation.

764 Regarding syllable selection, Paradigm 1 included six syllables: /ba/, /da/, /ga/, /bu/, /du/,  
765 and /gu/, each repeated 60 times in both the SA and SI tasks. Paradigm 2 included eight  
766 syllables: /ba/, /da/, /ga/, /pa/, /ta/, /ka/, /sa/, and /sha/, each repeated 45 times in both the SA  
767 and SI tasks.

768 Four methodological controls were implemented to eliminate potential articulatory muscle  
769 activation during the speech imagery task: (1) During preoperative training, subjects were  
770 required to complete a full iteration of the intraoperative paradigm and were explicitly  
771 instructed to avoid articulatory movements or voicing during SI tasks. (2) During preoperative  
772 training, a researcher (Z.Z.) closely monitored the subjects, providing immediate corrections  
773 to ensure the absence of visible articulatory movements and audible speech during the SI task.  
774 (3) During the intraoperative SI task, real-time video recordings of subjects' orofacial  
775 movements and corresponding audio recordings were monitored and acquired using the Brain  
776 Mapping Interactive Stimulation System (Shenzhen Sinorad Medical Electronics Co. Ltd)<sup>77</sup> to  
777 confirm the absence of motor and auditory output. (4) For the last four subjects (S6–S9),  
778 intraoperative electromyographic (EMG) recordings were obtained using needle electrodes  
779 placed on orofacial muscles (e.g., orbicularis oris, mylohyoid) to ensure that the included SI  
780 trials exhibited no detectable EMG activity exceeding the predefined threshold indicative of  
781 articulatory muscle engagement.

782

### 783 **ECoG data acquisition and preprocessing**

784 During the experimental tasks, neural signals were recorded from one or two 128-channel  
785 ECoG grids (8 × 16, 3 or 4 mm spacing) using a multichannel amplifier optically connected to  
786 a digital signal processor (Tucker-Davis Technologies). The local field potential at each  
787 electrode contact was amplified and sampled at 3052 Hz.

788 The raw voltage waveforms were visually inspected, and channels with undetectable signal  
789 variation relative to noise or exhibiting continuous epileptiform activity were excluded. Time  
790 segments containing electrical or movement-related artifacts were manually identified and  
791 removed. The remaining signals were then notch-filtered to eliminate line noise at 50 Hz, 100  
792 Hz, and 150 Hz.

793 To extract frequency-specific neural activity, band-pass filtering was applied to isolate the  
794 following frequency bands: theta (4–8 Hz), alpha (8–12 Hz), beta1 (12–24 Hz), beta2 (24–40  
795 Hz), low gamma (40–70 Hz), and high gamma (70–150 Hz)<sup>21,33</sup>. The analytic amplitude  
796 (envelope) of each band was computed using the Hilbert transform and smoothed using a  
797 Gaussian filter. The signal for each frequency band was obtained by averaging the analytic  
798 amplitude across eight logarithmically spaced sub-bands within the specified frequency range  
799<sup>33</sup>. Finally, the signal was down-sampled to 100 Hz and z-scored using the entire recording  
800 block for normalization<sup>31</sup>.

801

### 802 **Cortical surface extraction and electrode localization**

803 We used FreeSurfer to generate pial surface reconstructions from preoperative T1-weighted  
804 MRI scans. To localize electrode positions, the three-dimensional coordinates of the grid  
805 corners were recorded intraoperatively using the Medtronic neuronavigation system. These  
806 corner electrodes were then co-registered to the preoperative MRI, with intraoperative  
807 photographs serving as additional reference points. The remaining electrode positions were  
808 subsequently estimated using the “img\_pipe” package in Python, which interpolates and

809 extrapolates electrode locations based on the recorded corner coordinates <sup>78</sup>. Electrode  
810 positions in individual spaces were normalized to the MNI 152 (2009, asymmetric) template,  
811 and anatomical labels were assigned based on the Desikan-Killiany <sup>79</sup> and Human Connectome  
812 Project's multi-modal parcellation (MMP) atlases <sup>36</sup>.

813

### 814 **Responsive electrode selection**

815 After the neural activity computation in each frequency band, response electrodes for both SA  
816 and SI were identified. The average time difference between syllable onset in the audio (SA  
817 task) and the reference time point (ipsilateral key press in Paradigm 1 or the go-cue appearance  
818 in Paradigm 2) was calculated. This value was then used to align the ECoG data to the syllable  
819 onset in the SA modality or internal syllable onset in the SI modality (set as 0 ms), followed by  
820 segmentation and averaging across tasks.

821 The neural activity between -500 ms and -400 ms served as baseline resting-state activity.  
822 A two-sample *t*-test (Bonferroni correction,  $\alpha = 0.05/160$ , where 160 is the number of time  
823 points tested) compared average neural activity between baseline and each time point from -  
824 400 ms to +1200 ms. Electrode responses were considered significant if at least 100 ms (10  
825 consecutive time points) showed significant differences. To control for motor responses in  
826 Paradigm 1, electrodes responsive to the ipsilateral key press task (control) were excluded from  
827 further analysis.

828

### 829 **Speech feature extraction**

830 We employed a similar Praat parameter extraction approach in line with our previous research  
831 <sup>80,81</sup>. We extracted the pitch contour ( $\log F_0$ ,  $F_0$  = fundamental frequency) of each syllable from  
832 the audio recorded during the SA modality with an autocorrelation method in Praat  
833 (<https://www.fon.hum.uva.nl/praat/>, Version 6.1.01) <sup>82</sup>, which was subsequently used to  
834 compute the one-dimensional trajectory for the larynx. Additionally, we addressed halving and  
835 doubling errors during the extraction process. Individual pitch minimum and maximum values

836 were determined for each participant, and a timestep of 0.01s was employed. All other  
837 parameters adhered to Praat's default settings.

838 The onset (or internal onset for SI modality) of phonemes for all syllables was manually  
839 labeled in Praat by the researchers (Z.Z., Z.W., and Y.L.) and subsequently categorized into  
840 four places of articulation (POA) classes: (1) Labial: /b/, /p/; (2) Coronal: /d/, /t/, /s/, /sh/; (3)  
841 Dorsal: /g/, /k/; (4) Vocalic: /a/, /u/.

842

### 843 **Speaker-Independent Acoustic-to-Articulatory Inversion (AAI)**

844 Following prior studies<sup>28,83</sup>, we applied a Speaker-Independent Acoustic-to-Articulatory  
845 Inversion (AAI) algorithm to convert synchronized audio recorded during the SA modality into  
846 13-dimensional articulatory kinematic trajectories (AKTs) for each syllable. Previous studies  
847 have demonstrated a strong positive correlation between these articulatory kinematic  
848 trajectories and electromagnetic midsagittal articulography recordings<sup>28</sup>. These AKTs captured  
849 two-dimensional trajectories (x and y directions) for six articulators—upper lip (UL), lower lip  
850 (LL), lower incisor/jaw (LI), tongue tip (TT), tongue body (TB), and tongue dorsum (TD)—  
851 and a one-dimensional trajectory for the larynx (L), represented by the scaled log  $F_0$   
852 (fundamental frequency) and normalized to a range of -1 to 1 a.u. For SI tokens, we used the  
853 corresponding SA tokens' AKTs from the same trial as representations of internal AKTs. For  
854 detailed code of the Speaker-Independent AAI algorithm, please refer to:  
855 [github.com/articulatory/articulatory](https://github.com/articulatory/articulatory).

856

### 857 **Encoding model**

858 We used time-delayed linear encoding models, known as temporal receptive field models<sup>84</sup>, to  
859 evaluate what features drive the neural activity in frontoparietal regions during SA and SI  
860 modalities. Temporal receptive field (TRF) models predict neural activity using speech-related  
861 features (AKTs or POAs) in a time window (-400 ms to + 400 ms) around the neural activity.  
862 In particular, we fit the linear model for each electrode:

863 
$$y(t) = \sum_{f=1}^F \sum_{\tau=0}^T \boldsymbol{\beta}_f^T(\tau) \mathbf{x}_f(t - \tau) + \epsilon$$

864 Parameter  $y$  is the neural activity recorded from the electrode,  $\mathbf{x}_f(t - \tau)$  is the stimulus  
865 representation vector of feature set  $f$  at time  $t - \tau$ ,  $\boldsymbol{\beta}_f(\tau)$  is the regression weights for feature  
866 set  $f$  at time lag  $\tau$ , and  $\epsilon$  is the Gaussian noise.

867 To prevent overfitting, we employed L2 regularization and cross-validation. The data were  
868 divided into three mutually exclusive sets, comprising 80%, 10%, and 10% of the samples. The  
869 80% set was used for training, the 10% set for optimizing the L2 regularization hyperparameter,  
870 and the remaining 10% set for testing. Model performance was evaluated by the coefficient of  
871 determination ( $R^2$ ) between the actual and predicted neural activity values on the held-out test  
872 data. This procedure was repeated using five-fold cross-validation, with the final model  
873 performance reported as the mean  $R^2$  across all folds, reflecting the total variance explained  
874 by the model.

875 In the full TRF model for AKTs, the feature matrix of the 13-dimensional trajectories for  
876 seven articulators was included. To compare the performance of the AKT models across all  
877 frequency bands in both the SA and SI modalities, we first conducted a Kruskal-Wallis  $H$ -test  
878 with a significance threshold of  $p < 0.05$ . Subsequently, pairwise Mann-Whitney  $U$  tests were  
879 performed, with a significance threshold of  $p < 0.05$  after the Benjamini-Hochberg correction.  
880 To quantify the unique contribution of each articulator, we fitted TRF models excluding each  
881 articulator in turn and computed the change in  $R^2$  ( $\Delta R^2$ ) between the full and reduced models.

882 We also trained a full TRF model incorporating POA cluster-based features, where the  
883 presence or absence of each phoneme's POA feature at its onset was marked using binary  
884 indicators (1 for presence, 0 for absence). The unique contribution of each POA was also  
885 quantified as the  $\Delta R^2$  between the full and reduced models.

886 Building upon the electrodes exhibiting significant activity in the beta1 band, we further  
887 excluded electrodes that did not encode AKTs. Specifically, we removed electrodes with an  $R^2$   
888 value of  $\leq 0.01$  in both the SA- and SI-AKT modalities, corresponding to the 50% threshold  
889 of the  $R^2$  distribution across all responsive electrodes (**Extended Data Fig. 4b**).

890

### 891 **K-means clustering for encoding electrode classification**

892 For dual-responsive electrodes, we aimed to determine whether beta1 neural activity encodes  
893 AKTs similarly across SA and SI modalities. To address this, we first calculated the  $R^2$   
894 difference between the SA-AKT and SI-AKT models for each dual-responsive electrode. We  
895 applied the Duda-Hart test to evaluate whether the differences surpassed the threshold (Duda-  
896 Hart statistic  $d > 1.645$ ,  $p < 0.05$ ), which would indicate the necessity for classification into at  
897 least two clusters<sup>85</sup>. Following this, we performed unsupervised K-means clustering (KMeans  
898 from scikit-learn python package), iterating the number of clusters from 2 to 10, and selected  
899 the optimal number of clusters based on the highest silhouette score<sup>86</sup>. For each cluster, we  
900 analyzed the similarity in performance ( $R^2$  values) between the SA-AKT and SI-AKT models  
901 for the electrodes within the cluster using Pearson correlation analysis. Differences in  
902 performance were assessed using the Wilcoxon signed-rank test (two-sided). A  $p$ -value of  $<$   
903 0.05 was considered statistically significant for both analyses.

904

### 905 **Cosine similarity of encoding patterns across SA and SI modalities**

906 To assess whether electrodes within specific clusters exhibit shared AKT encoding patterns, we  
907 constructed encoding matrices for each modality. Each matrix entry represented each  
908 articulator's unique contribution ( $\Delta R^2$ ) across electrodes, scaled to a range of 0 to 1 arbitrary  
909 unit. We then calculated the cosine similarity between the encoding matrices for the two  
910 modalities. A permutation test was performed to evaluate statistical significance: electrode  
911 labels within each encoding matrix were randomly shuffled, and the cosine similarity was  
912 recalculated for 10,000 iterations to generate a null distribution. The observed cosine similarity  
913 was compared to this distribution, with significance determined if the  $p$ -value was  $< 0.05$ . The  
914 same approach was applied to assess the cosine similarity of POA encoding patterns between  
915 the SA and SI modalities.

916

## 917 **Hierarchical clustering**

918 Agglomerative hierarchical clustering was performed using Ward's method. Clustering was  
919 applied along the electrode dimension for AKT (**Fig. 2s, 2u**) and POA encoding matrices  
920 (**Extended Data Fig. 5a, 5c**).

921 To examine clustering patterns of AKTs across phonemes, we constructed a phoneme-  
922 specific AKT displacement matrix (**Fig. 4a**). For each phoneme, we calculated the average  
923 AKTs across all subjects, defining the direction of maximum deviation from the resting state  
924 as the primary movement direction (positive). The maximum displacement of each phoneme's  
925 13 AKTs was then computed, with a displacement sign assigned based on alignment with the  
926 primary movement direction. The resulting AKT displacement matrix for each phoneme was  
927 normalized to a range of 0 to 1 arbitrary unit. Hierarchical clustering was performed on both  
928 the phoneme and AKT dimensions using SciPy and Seaborn Python packages <sup>87</sup>.

929

## 930 **Spatial distribution of encoding electrodes**

931 To investigate the spatial distribution of encoding electrodes across different categories (SA-  
932 SI shared, SA-specific, and SI-specific), kernel density estimation (KDE) was performed in the  
933 y-z plane of MNI space. KDE was applied with Scott's bandwidth selection method to construct  
934 probability density distributions for each encoding category. The distributions were thresholded  
935 at 0.5 cumulative density (a.u.).

936 The KDE maps were processed using a maximum filter from the SciPy package  
937 (`scipy.ndimage.maximum_filter`) to identify local density peaks. Statistical comparisons of  
938 electrode distributions in MNI space were performed using the two-sided Kolmogorov-  
939 Smirnov test to assess differences along the y-axis and z-axis (significant if  $p < 0.05$ ).

940 Similar KDE maps were generated for each POA across different encoding categories, with  
941 electrode positions weighted by their unique contribution ( $\Delta R^2$ ) to the full POA model. To  
942 evaluate differences in weighted KDE distributions between POAs along the y-axis and z-axis,  
943 we applied a permutation test based on the Wasserstein distance. For each permutation (5000

944 iterations), electrode positions and weights ( $\Delta R^2$ ) were shuffled, and the Wasserstein distance  
945 was recalculated between POA distributions. The  $p$ -values were derived from the permutation  
946 distribution and corrected for multiple comparisons using the Benjamini-Hochberg procedure,  
947 with significance at  $p < 0.05$  after correction.

948

### 949 **Temporal distribution of encoding electrodes**

950 To investigate the temporal distribution of encoding electrodes, we calculated the onset and  
951 peak times during syllable production and imagery tasks. For each time point within the -400  
952 ms to +1200 ms window, an independent two-sample t-test was performed comparing task-  
953 related beta1 activity to baseline. The onset time was defined as the first time point with a  
954 sustained  $p$ -value  $< 0.01$  for at least 100 ms. The peak time was determined as the point of  
955 maximal beta1 amplitude following onset. Differences in onset and peak times across encoding  
956 categories were evaluated using the Mann-Whitney  $U$  test, with significance at  $p < 0.05$ .

957

### 958 **Triple-stream parallel neural decoding framework**

959 To determine whether beta1 activity across discrete multi-electrode sites could be leveraged  
960 for speech decoding, particularly for inner speech decoding in the SI modality, we developed  
961 a triple-stream parallel neural decoding framework (**Fig. 5a**) using the Pytorch and Keras  
962 packages<sup>88,89</sup>. This framework decodes beta1 activity from SA- and SI-responsive electrodes  
963 into three highly interrelated components: syllable classification (syllable classifier, **Fig. 5b**),  
964 audio synthesis (speech synthesizer, **Fig. 5c**), and AKT reconstruction (articulatory movement  
965 synthesizer, **Fig. 5d**). Notably, since the SI modality does not generate auditory or motor output,  
966 we incorporated the speech feature matrix from the corresponding SA trial as an internal  
967 representation of speech to train the decoding model. For all three streams, input data were  
968 structured as channels  $\times$  time steps, where channels corresponded to SA- or SI-responsive  
969 electrodes, depending on the task modality. The time window extended from 0.5 s before to 1.5  
970 s after syllable onset, comprising 200 time steps at a sampling rate of 100 Hz. For the  
971 articulatory movement synthesizer, the first and last four time steps were excluded to mitigate

972 the edge effects of AKTs.

973 **The syllable classifier (Fig. 5b)** is an adapted 34-layer Residual Network (ResNet-34)  
974 architecture<sup>90</sup>. The first convolutional layer employed a  $7 \times 7$  kernel, followed by a batch  
975 normalization layer and a 2D max-pooling layer. Four residual blocks (ResBlocks) were  
976 subsequently applied. All convolutional layers in the decoder utilized Rectified Linear Units  
977 (ReLU) as activation functions<sup>91</sup>. A 40% dropout layer was included to mitigate overfitting.  
978 The network was trained to minimize joint cross-entropy loss using the stochastic gradient  
979 descent (SGD) optimizer, with an initial learning rate of 0.01, a decay rate of  $1 \times 10^{-10^6}$ , the  
980 momentum of 0.9, and Nesterov acceleration enabled. Optimization was halted when  
981 validation loss ceased decreasing. The classifier's performance was evaluated using syllable  
982 decoding accuracy, employing a 10-fold cross-validation scheme to determine the average  
983 performance across folds.

984 **The speech synthesizer (Fig. 5c)** employs a deep convolutional neural network  
985 architecture composed of four sequential modules, each consisting of a Conv2D layer, a  
986 ResBlock<sup>90</sup>, and a 2D max-pooling layer, followed by a fully connected layer. All  
987 convolutional kernels were set to  $3 \times 3$ . To assess the quality of synthesized speech waveforms,  
988 we used Mel-cepstral distortion (MCD) as an objective measure. MCD quantifies the error in  
989 Mel-Frequency Cepstral Coefficients (MFCCs)<sup>30,92</sup> and was calculated as follows:

$$990 \quad MCD \text{ (dB)} = \frac{10}{\ln 10} \sqrt{2 \sum_{d=1}^{24} (mc_d^{(\hat{y})} - mc_d^{(y)})^2}$$

991 where  $d$  represents each MFCC dimension ( $0 < d < 25$ ),  $\hat{y}$  is the synthesized speech, and  
992  $y$  is the actual participant-produced acoustic signal<sup>92</sup>. Performance evaluation used a 10-fold  
993 cross-validation strategy to determine the average MCD across folds.

994 **The articulatory movement (Fig. 5d)** synthesizer employs a vector-quantized variational  
995 autoencoder (VQ-VAE)<sup>93</sup>. In the first stage, an encoder block composed of residual and non-  
996 local blocks extracts features from the ECoG signals. In the second stage, extracted features  
997 are mapped onto a 256-dimensional articulatory codebook, yielding a discretized codebook  
998 representation. In the third stage, the codebook mapping is fed into a decoder block (also

999 consisting of residual and non-local blocks)<sup>90</sup> to reconstruct articulatory movements with  
1000 dimensions of  $13 \times 192$ . The network was optimized to minimize mean squared error loss using  
1001 the Adam optimizer, with an initial learning rate of  $2.25 \times 10^{-5}$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.9$ , and  $\varepsilon = 1 \times 10^{-8}$ .  
1002 Performance evaluation was conducted using a 5-fold cross-validation strategy. The dataset  
1003 was randomly divided into ten subsets, with eight used for training, one for validation, and one  
1004 for testing. For each fold, performance was assessed based on the mean Pearson's correlation  
1005 coefficient ( $r$ ) across 13-dimensional AKTs. The final performance was reported as the average  
1006 across all five folds.

1007 To validate the robustness of this triple-stream parallel decoding framework, we conducted  
1008 tests on both SA and SI modalities in nine participants. Median syllable decoding accuracy was  
1009 compared against chance levels (1/6 for six-class tasks in subjects S1-S5; 1/8 for eight-class  
1010 tasks in subjects S6-S9). Median MCD was compared against the 8 dB threshold, the upper  
1011 limit considered acceptable for voice recognition systems<sup>30</sup>. Additionally, median  $r$  values  
1012 were compared against chance levels, defined as the correlation between AKTs synthesized by  
1013 an untrained, randomly initialized model and actual AKTs. All statistical comparisons were  
1014 conducted using the Mann-Whitney  $U$  test. Modality differences were assessed using the  
1015 Wilcoxon signed-rank test, with  $p < 0.05$  considered statistically significant.

1016 To further investigate whether cross-modal articulatory movement decoding could be  
1017 achieved using SA-SI shared electrodes, we trained four articulatory movement synthesizers:  
1018 SA2SA, SA2SI, SI2SI, and SI2SA (naming format indicates "training set to testing set," e.g.,  
1019 SA2SI denotes training with SA modality beta1 activity and testing with SI modality beta1  
1020 activity to evaluate the correlation between synthesized and ground-truth movements).  
1021 Differences among these four models were assessed using the Kruskal-Wallis  $H$  test, with  $p <$   
1022  $0.05$  considered statistically significant.

1023 Lastly, we examined whether the articulatory movement synthesizer effectively captured  
1024 a stable relationship between beta1 activity and AKTs, allowing for extrapolative  
1025 generalization. Specifically, the model's ability to predict movement trajectories of syllables  
1026 excluded from the training set was assessed using a Leave-One-Syllable-Out strategy. In each  
1027 iteration, one syllable was excluded from training, and the model was trained on the remaining

1028 syllables. The excluded syllable was then used as a testing set, and the correlation ( $r$ ) between  
1029 predicted and actual AKTs was computed. Chance-level  $r$  values were determined using an  
1030 untrained, randomly initialized model. Median  $r$  values and modality differences were  
1031 statistically assessed using the Mann-Whitney  $U$  and Wilcoxon signed-rank tests, respectively.

1032

### 1033 **Data availability**

1034 The data set generated during the current study will be made available from the authors upon  
1035 reasonable request.

1036

### 1037 **Code availability**

1038 The completely developed code that operates on the full data set will be made available from  
1039 the corresponding authors upon reasonable request.

1040

### 1041 **Acknowledgments**

1042 Dr. Junfeng Lu is supported by STI 2030—Major Projects (2022ZD0212300) and National  
1043 Natural Science Foundation of China (32371146). Dr. Zehao Zhao is supported by the  
1044 Postdoctoral Fellowship Program of CPSF (GZB20240661). Dr. Jinsong Wu is supported by  
1045 Innovation Program of Shanghai Municipal Education Commission (2023ZKZD13). Dr.  
1046 Yuanning Li is supported by the National Natural Science Foundation of China (32371154),  
1047 Shanghai Rising-Star Program (24QA2705500), Shanghai Pujiang Program (22PJ1410500,  
1048 Y.L.), and the Lingang Laboratory (LG-GG-202402-06). The computations in this research are  
1049 supported by the HPC Platform of ShanghaiTech University.

1050

### 1051 **Author contributions**

1052 Z.Z., Z.W., Y.Liu., Y.Li., J.L., and J.W. conceived and designed the study. Z.Z., Z.W., Y.Li.,  
1053 and Y.Y. constructed the temporal receptive field model and performed statistical analyses. Z.Z.,  
1054 Z.W., Y.Li., and Y.Liu. developed and optimized the decoding framework. Z.Z., Y.Liu., Y.Q.,  
1055 J.L., and J.W. collected neural and behavioral data. J.L. and J.W. coordinated session logistics,  
1056 managed equipment, and performed cortical electrode placement. X.G., B.Y., S.X.T., and X.T.  
1057 provided guidance on neurolinguistic experimental design and contributed to manuscript  
1058 revision. J.W. supervised the overall project and provided strategic guidance. The study was  
1059 jointly supervised by G.C., Y.Li., J.L., and J.W. Z.Z. wrote the initial draft with input from  
1060 Y.Liu. and Z.W. All authors reviewed and edited the manuscript.

1061

## 1062 **Competing interests**

1063 The authors report no competing interests.

1064

## 1065 **References**

- 1066 1 Moulton, S. T. & Kosslyn, S. M. Imagining predictions: mental imagery as mental emulation.  
1067 *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **364**, 1273-1280,  
1068 doi:10.1098/rstb.2008.0314 (2009).
- 1069 2 Kosslyn, S. M., Ganis, G. & Thompson, W. L. Neural foundations of imagery. *Nat Rev Neurosci* **2**, 635-  
1070 642, doi:10.1038/35090055 (2001).
- 1071 3 Pearson, J., Naselaris, T., Holmes, E. A. & Kosslyn, S. M. Mental Imagery: Functional Mechanisms and  
1072 Clinical Applications. *Trends Cogn Sci* **19**, 590-602, doi:10.1016/j.tics.2015.08.003 (2015).
- 1073 4 Schacter, D. L., Addis, D. R. & Buckner, R. L. Remembering the past to imagine the future: the  
1074 prospective brain. *Nat Rev Neurosci* **8**, 657-661, doi:10.1038/nrn2213 (2007).
- 1075 5 Schacter, D. L. *et al.* The future of memory: remembering, imagining, and the brain. *Neuron* **76**, 677-  
1076 694, doi:10.1016/j.neuron.2012.11.001 (2012).
- 1077 6 Grandchamp, R. *et al.* The ConDialInt Model: Condensation, Dialogality, and Intentionality Dimensions  
1078 of Inner Speech Within a Hierarchical Predictive Control Framework. *Front Psychol* **10**, 2019,  
1079 doi:10.3389/fpsyg.2019.02019 (2019).
- 1080 7 Kearney, E. & Guenther, F. H. Articulating: The Neural Mechanisms of Speech Production. *Lang Cogn*  
1081 *Neurosci* **34**, 1214-1229, doi:10.1080/23273798.2019.1589541 (2019).
- 1082 8 Buzsaki, G. & Tingley, D. Cognition from the Body-Brain Partnership: Exaptation of Memory. *Annu Rev*  
1083 *Neurosci* **46**, 191-210, doi:10.1146/annurev-neuro-101222-110632 (2023).
- 1084 9 Duffau, H., Moritz-Gasser, S. & Mandonnet, E. A re-examination of neural basis of language processing:

- 1085 proposal of a dynamic hodotopical model from data provided by brain stimulation mapping during  
1086 picture naming. *Brain Lang* **131**, 1-10, doi:10.1016/j.bandl.2013.05.011 (2014).
- 1087 10 Hickok, G., Venezia, J. & Teghipco, A. Beyond broca: neural architecture and evolution of a dual motor  
1088 speech coordination system. *PsyArXiv [Preprints]* doi **10** (2021).
- 1089 11 Cooney, C., Folli, R. & Coyle, D. Neurolinguistics Research Advancing Development of a Direct-Speech  
1090 Brain-Computer Interface. *iScience* **8**, 103-125, doi:10.1016/j.isci.2018.09.016 (2018).
- 1091 12 Alderson-Day, B. & Fernyhough, C. Inner Speech: Development, Cognitive Functions, Phenomenology,  
1092 and Neurobiology. *Psychol Bull* **141**, 931-965, doi:10.1037/bul0000021 (2015).
- 1093 13 Rodriguez-Fornells, A., León-Cabrera, P., Gabarras, A. & Sierpowska, J. in *Intraoperative Mapping of*  
1094 *Cognitive Networks* Ch. Chapter 23, 381-409 (2021).
- 1095 14 Oppenheim, G. M. & Dell, G. S. Motor movement matters: the flexible abstractness of inner speech.  
1096 *Mem Cognit* **38**, 1147-1160, doi:10.3758/MC.38.8.1147 (2010).
- 1097 15 Wheeldon, L. R. & Levelt, W. J. Monitoring the time course of phonological encoding. *Journal of*  
1098 *memory and language* **34**, 311-334 (1995).
- 1099 16 Tian, X. & Poeppel, D. Mental imagery of speech and movement implicates the dynamics of internal  
1100 forward models. *Front Psychol* **1**, 166, doi:10.3389/fpsyg.2010.00166 (2010).
- 1101 17 Tian, X. & Poeppel, D. Mental imagery of speech: linking motor and perceptual systems through internal  
1102 simulation and estimation. *Front Hum Neurosci* **6**, doi:10.3389/fnhum.2012.00314 (2012).
- 1103 18 Tian, X., Zarate, J. M. & Poeppel, D. Mental imagery of speech implicates two mechanisms of perceptual  
1104 reactivation. *Cortex* **77**, 1-12, doi:10.1016/j.cortex.2016.01.002 (2016).
- 1105 19 Pei, X. *et al.* Spatiotemporal dynamics of electrocorticographic high gamma activity during overt and  
1106 covert word repetition. *NeuroImage* **54**, 2960-2972, doi:10.1016/j.neuroimage.2010.10.029 (2011).
- 1107 20 Proix, T. *et al.* Imagined speech can be decoded from low- and cross-frequency intracranial EEG features.  
1108 *Nat Commun* **13**, 48, doi:10.1038/s41467-021-27725-3 (2022).
- 1109 21 de Borman, A. *et al.* Imagined speech event detection from electrocorticography and its transfer between  
1110 speech modes and subjects. *Commun Biol* **7**, 818, doi:10.1038/s42003-024-06518-6 (2024).
- 1111 22 Goutman, S. A. *et al.* Recent advances in the diagnosis and prognosis of amyotrophic lateral sclerosis.  
1112 *The Lancet. Neurology* **21**, 480-493, doi:10.1016/S1474-4422(21)00465-8 (2022).
- 1113 23 Anumanchipalli, G. K., Chartier, J. & Chang, E. F. Speech synthesis from neural decoding of spoken  
1114 sentences. *Nature* **568**, 493-498, doi:10.1038/s41586-019-1119-1 (2019).
- 1115 24 Moses, D. A. *et al.* Neuroprosthesis for Decoding Speech in a Paralyzed Person with Anarthria. *N Engl*  
1116 *J Med* **385**, 217-227, doi:10.1056/NEJMoa2027540 (2021).
- 1117 25 Metzger, S. L. *et al.* A high-performance neuroprosthesis for speech decoding and avatar control. *Nature*  
1118 **620**, 1037-1046, doi:10.1038/s41586-023-06443-4 (2023).
- 1119 26 Willett, F. R. *et al.* A high-performance speech neuroprosthesis. *Nature* **620**, 1031-1036,  
1120 doi:10.1038/s41586-023-06377-x (2023).
- 1121 27 Bouchard, K. E., Mesgarani, N., Johnson, K. & Chang, E. F. Functional organization of human  
1122 sensorimotor cortex for speech articulation. *Nature* **495**, 327-332, doi:10.1038/nature11911 (2013).
- 1123 28 Chartier, J., Anumanchipalli, G. K., Johnson, K. & Chang, E. F. Encoding of Articulatory Kinematic  
1124 Trajectories in Human Speech Sensorimotor Cortex. *Neuron* **98**, 1042-1054 e1044,  
1125 doi:10.1016/j.neuron.2018.04.031 (2018).
- 1126 29 Dichter, B. K., Breshears, J. D., Leonard, M. K. & Chang, E. F. The Control of Vocal Pitch in Human  
1127 Laryngeal Motor Cortex. *Cell* **174**, 21-+, doi:10.1016/j.cell.2018.05.016 (2018).
- 1128 30 Liu, Y. *et al.* Decoding and synthesizing tonal language speech from brain activity. *Science Advances* **9**,

- 1129 eadh0478, doi:10.1126/sciadv.adh0478 (2023).
- 1130 31 Lu, J. *et al.* Neural control of lexical tone production in human laryngeal motor cortex. *Nat Commun* **14**,  
1131 6917, doi:10.1038/s41467-023-42175-9 (2023).
- 1132 32 Smith, E. & Delargy, M. Locked-in syndrome. *Bmj* **330**, 406-409 (2005).
- 1133 33 Hamilton, L. S., Oganian, Y., Hall, J. & Chang, E. F. Parallel and distributed encoding of speech across  
1134 human auditory cortex. *Cell* **184**, 4626-4639.e4613, doi:10.1016/j.cell.2021.07.019 (2021).
- 1135 34 Theunissen, F. E., Sen, K. & Doupe, A. J. Spectral-temporal receptive fields of nonlinear auditory  
1136 neurons obtained using natural sounds. *J Neurosci* **20**, 2315-2331 (2000).
- 1137 35 Linden, H. *et al.* Modeling the spatial reach of the LFP. *Neuron* **72**, 859-872,  
1138 doi:10.1016/j.neuron.2011.11.006 (2011).
- 1139 36 Glasser, M. F. *et al.* A multi-modal parcellation of human cerebral cortex. *Nature* **536**, 171-178,  
1140 doi:10.1038/nature18933 (2016).
- 1141 37 Willett, F. R. *et al.* Hand Knob Area of Premotor Cortex Represents the Whole Body in a Compositional  
1142 Way. *Cell* **181**, 396-409.e326, doi:10.1016/j.cell.2020.02.043 (2020).
- 1143 38 Lu, J. *et al.* Functional maps of direct electrical stimulation-induced speech arrest and anomia: a  
1144 multicentre retrospective study. *Brain* **144**, 2541-2553, doi:10.1093/brain/awab125 (2021).
- 1145 39 Gordon, E. M. *et al.* A somato-cognitive action network alternates with effector regions in motor cortex.  
1146 *Nature* **617**, 351-359, doi:10.1038/s41586-023-05964-2 (2023).
- 1147 40 Natraj, N. *et al.* Sampling representational plasticity of simple imagined movements across days enables  
1148 long-term neuroprosthetic control. *Cell* **188**, 1208-1225.e1232, doi:10.1016/j.cell.2025.02.001 (2025).
- 1149 41 Morgan, A. M. *et al.* A magnitude-independent neural code for linguistic information during sentence  
1150 production. *bioRxiv*, 2024.2006.2020.599931 (2025).
- 1151 42 Kilavik, B. E., Zaepffel, M., Brovelli, A., MacKay, W. A. & Riehle, A. The ups and downs of beta  
1152 oscillations in sensorimotor cortex. *Exp Neurol* **245**, 15-26, doi:10.1016/j.expneurol.2012.09.014 (2013).
- 1153 43 Saleh, M., Reimer, J., Penn, R., Ojakangas, C. L. & Hatsopoulos, N. G. Fast and slow oscillations in  
1154 human primary motor cortex predict oncoming behaviorally relevant cues. *Neuron* **65**, 461-471,  
1155 doi:10.1016/j.neuron.2010.02.001 (2010).
- 1156 44 Gilbertson, T. *et al.* Existing motor state is favored at the expense of new movement during 13-35 Hz  
1157 oscillatory synchrony in the human corticospinal system. *J Neurosci* **25**, 7771-7779,  
1158 doi:10.1523/jneurosci.1762-05.2005 (2005).
- 1159 45 Rubino, D., Robbins, K. A. & Hatsopoulos, N. G. Propagating waves mediate information transfer in the  
1160 motor cortex. *Nat Neurosci* **9**, 1549-1557, doi:10.1038/nn1802 (2006).
- 1161 46 Sanes, J. N. & Donoghue, J. P. Oscillations in local field potentials of the primate motor cortex during  
1162 voluntary movement. *Proceedings of the National Academy of Sciences of the United States of America*  
1163 **90**, 4470-4474, doi:10.1073/pnas.90.10.4470 (1993).
- 1164 47 Crone, N. E. *et al.* Functional mapping of human sensorimotor cortex with electrocorticographic spectral  
1165 analysis. I. Alpha and beta event-related desynchronization. *Brain* **121** ( Pt 12), 2271-2299,  
1166 doi:10.1093/brain/121.12.2271 (1998).
- 1167 48 Penfield, W. & Boldrey, E. Somatic motor and sensory representation in the cerebral cortex of man as  
1168 studied by electrical stimulation. *Brain* **60**, 389-443, doi:DOI 10.1093/brain/60.4.389 (1937).
- 1169 49 Zhao, Z. *et al.* Convergence of the arcuate fasciculus and third branch of the superior longitudinal  
1170 fasciculus with direct cortical stimulation-induced speech arrest area in the anterior ventral precentral  
1171 gyrus. *J Neurosurg* **139**, 1140-1151, doi:10.3171/2023.1.JNS222575 (2023).
- 1172 50 Neta, M. *et al.* Spatial and temporal characteristics of error-related activity in the human brain. *J Neurosci*

- 1173           **35**, 253-266, doi:10.1523/JNEUROSCI.1313-14.2015 (2015).
- 1174    51       Dum, R. P., Levinthal, D. J. & Strick, P. L. Motor, cognitive, and affective areas of the cerebral cortex  
1175           influence the adrenal medulla. *Proceedings of the National Academy of Sciences of the United States of*  
1176           *America* **113**, 9922-9927, doi:10.1073/pnas.1605044113 (2016).
- 1177    52       Duffau, H. Stimulation mapping of white matter tracts to study brain functional connectivity. *Nat Rev*  
1178           *Neurol* **11**, 255-265, doi:10.1038/nrneurol.2015.51 (2015).
- 1179    53       Yeh, F.-C. Population-based tract-to-region connectome of the human brain and its hierarchical topology.  
1180           *Nature Communications* **13**, 4933, doi:10.1038/s41467-022-32595-4 (2022).
- 1181    54       Zhao, Z., Liu, Y., Zhang, J., Lu, J. & Wu, J. Where is the speech production area? Evidence from direct  
1182           cortical electrical stimulation mapping. *Brain* **144**, e61, doi:10.1093/brain/awab178 (2021).
- 1183    55       Maldonado, I. L., Moritz-Gasser, S. & Duffau, H. Does the left superior longitudinal fascicle subservice  
1184           language semantics? A brain electrostimulation study. *Brain Struct Funct* **216**, 263-274,  
1185           doi:10.1007/s00429-011-0309-x (2011).
- 1186    56       Sarubbo, S. *et al.* Mapping critical cortical hubs and white matter pathways by direct electrical  
1187           stimulation: an original functional atlas of the human brain. *NeuroImage* **205**, 116237,  
1188           doi:10.1016/j.neuroimage.2019.116237 (2020).
- 1189    57       Castellucci, G. A., Kovach, C. K., Howard, M. A., 3rd, Greenlee, J. D. W. & Long, M. A. A speech  
1190           planning network for interactive language use. *Nature* **602**, 117-122, doi:10.1038/s41586-021-04270-z  
1191           (2022).
- 1192    58       Kunz, E. M. *et al.* Representation of verbal thought in motor cortex and implications for speech  
1193           neuroprostheses. *bioRxiv*, 2024.2010.2004.616375 (2024).
- 1194    59       Wandelt, S. K. *et al.* Representation of internal speech by single neurons in human supramarginal gyrus.  
1195           *Nat Hum Behav* **8**, 1136-1149, doi:10.1038/s41562-024-01867-y (2024).
- 1196    60       Belyk, M. & Brown, S. The origins of the vocal brain in humans. *Neurosci Biobehav Rev* **77**, 177-193,  
1197           doi:10.1016/j.neubiorev.2017.03.014 (2017).
- 1198    61       Roux, F. E., Niare, M., Charni, S., Giussani, C. & Durand, J. B. Functional architecture of the motor  
1199           homunculus detected by electrostimulation. *J Physiol* **598**, 5487-5504, doi:10.1113/JP280156 (2020).
- 1200    62       Zhou, Y. *et al.* Electrical stimulation-induced speech-related negative motor responses in the lateral  
1201           frontal cortex. *J Neurosurg* **137**, 496-504, doi:10.3171/2021.9.JNS211069 (2022).
- 1202    63       Liang, B., Li, Y., Zhao, W. & Du, Y. Bilateral human laryngeal motor cortex in perceptual decision of  
1203           lexical tone and voicing of consonant. *Nat Commun* **14**, 4710, doi:10.1038/s41467-023-40445-0 (2023).
- 1204    64       Indefrey, P. The spatial and temporal signatures of word production components: a critical update. *Front*  
1205           *Psychol* **2**, 255, doi:10.3389/fpsyg.2011.00255 (2011).
- 1206    65       Indefrey, P. & Levelt, W. J. The spatial and temporal signatures of word production components.  
1207           *Cognition* **92**, 101-144 (2004).
- 1208    66       Levelt, W. J., Praamstra, P., Meyer, A. S., Helenius, P. & Salmelin, R. An MEG study of picture naming.  
1209           *J Cognitive Neurosci* **10**, 553-567 (1998).
- 1210    67       Solodkin, A., Hlustik, P., Chen, E. E. & Small, S. L. Fine modulation in network activation during motor  
1211           execution and motor imagery. *Cerebral cortex* **14**, 1246-1255, doi:10.1093/cercor/bhh086 (2004).
- 1212    68       Perrone-Bertolotti, M., Rapin, L., Lachaux, J. P., Baciú, M. & Loevenbruck, H. What is that little voice  
1213           inside my head? Inner speech phenomenology, its role in cognitive performance, and its relation to self-  
1214           monitoring. *Behav Brain Res* **261**, 220-239, doi:10.1016/j.bbr.2013.12.034 (2014).
- 1215    69       Hickok, G. The architecture of speech production and the role of the phoneme in speech processing.  
1216           *Language, Cognition and Neuroscience* **29**, 2-20 (2014).

- 1217 70 Hickok, G. & Poeppel, D. Dorsal and ventral streams: a framework for understanding aspects of the  
1218 functional anatomy of language. *Cognition* **92**, 67-99, doi:10.1016/j.cognition.2003.10.011 (2004).
- 1219 71 Silva, A. B. *et al.* A Neurosurgical Functional Dissection of the Middle Precentral Gyrus during Speech  
1220 Production. *J Neurosci* **42**, 8416-8426, doi:10.1523/JNEUROSCI.1614-22.2022 (2022).
- 1221 72 Amunts, K. *et al.* Broca's region: novel organizational principles and multiple receptor mapping. *PLoS*  
1222 *Biol* **8**, e1000489 (2010).
- 1223 73 Judaš, M. & Ceganec, M. Oskar Vogt: The first myeloarchitectonic map of the human frontal cortex.  
1224 *Translational Neuroscience* **1**, 72-94, doi:10.2478/v10134-010-0005-z (2010).
- 1225 74 Garey, L. J. *Brodman's' localisation in the cerebral cortex'*. (World Scientific, 1999).
- 1226 75 Chen, X. *et al.* A neural speech decoding framework leveraging deep learning and speech synthesis.  
1227 *Nature Machine Intelligence* **6**, 467-480, doi:10.1038/s42256-024-00824-8 (2024).
- 1228 76 Littlejohn, K. T. *et al.* A streaming brain-to-voice neuroprosthesis to restore naturalistic communication.  
1229 *Nat Neurosci*, doi:10.1038/s41593-025-01905-6 (2025).
- 1230 77 Hameed, N. U. F. *et al.* A Novel Intraoperative Brain Mapping Integrated Task-Presentation Platform.  
1231 *Oper Neurosurg (Hagerstown)* **20**, 477-483, doi:10.1093/ons/opaa476 (2021).
- 1232 78 Hamilton, L. S., Chang, D. L., Lee, M. B. & Chang, E. F. Semi-automated Anatomical Labeling and  
1233 Inter-subject Warping of High-Density Intracranial Recording Electrodes in Electrocorticography. *Front*  
1234 *Neuroinform* **11**, 62, doi:10.3389/fninf.2017.00062 (2017).
- 1235 79 Desikan, R. S. *et al.* An automated labeling system for subdividing the human cerebral cortex on MRI  
1236 scans into gyral based regions of interest. *NeuroImage* **31**, 968-980,  
1237 doi:10.1016/j.neuroimage.2006.01.021 (2006).
- 1238 80 Tang, C., Hamilton, L. S. & Chang, E. F. Intonational speech prosody encoding in the human auditory  
1239 cortex. *Science* **357**, 797-801, doi:10.1126/science.aam8577 (2017).
- 1240 81 Li, Y., Tang, C., Lu, J., Wu, J. & Chang, E. F. Human cortical encoding of pitch in tonal and non-tonal  
1241 languages. *Nat Commun* **12**, 1161, doi:10.1038/s41467-021-21430-x (2021).
- 1242 82 Boersma, P. in *Proceedings of the institute of phonetic sciences*. 97-110 (Amsterdam).
- 1243 83 Wu, P. *et al.* in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal*  
1244 *Processing (ICASSP)*. 1-5.
- 1245 84 Theunissen, F. E. *et al.* Estimating spatio-temporal receptive fields of auditory and visual neurons from  
1246 their responses to natural stimuli. *Network: Computation in Neural Systems* **12**, 289 (2001).
- 1247 85 Murtagh, F. & Farid, M. M. Pattern Classification, by Richard O. Duda, Peter E. Hart, and David G.  
1248 Stork. **18**, 273-275 (2001).
- 1249 86 Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis.  
1250 *Journal of computational and applied mathematics* **20**, 53-65 (1987).
- 1251 87 Müllner, D. Modern hierarchical, agglomerative clustering algorithms. *arXiv preprint arXiv:1109.2378*  
1252 (2011).
- 1253 88 Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library. *Advances in*  
1254 *neural information processing systems* **32** (2019).
- 1255 89 Gulli, A. & Pal, S. *Deep learning with Keras*. (Packt Publishing Ltd, 2017).
- 1256 90 He, K., Zhang, X., Ren, S. & Sun, J. in *Proceedings of the IEEE conference on computer vision and*  
1257 *pattern recognition*. 770-778.
- 1258 91 Nair, V. & Hinton, G. E. in *Proceedings of the 27th international conference on machine learning (ICML-*  
1259 *10)*. 807-814.
- 1260 92 Kubichek, R. in *Proceedings of IEEE pacific rim conference on communications computers and signal*

- 1261            *processing*. 125-128 (IEEE).
- 1262    93        Esser, P., Rombach, R. & Ommer, B. in *Proceedings of the IEEE/CVF conference on computer vision*
- 1263            *and pattern recognition*. 12873-12883.
- 1264