1
2  *Transformation of a temporal speech cue to a spatial neural code in human auditory cortex*
3
4  Neal P. Fox[a], Matthew K. Leonard[a], Matthias J. Sjerps[b,c] & Edward F. Chang[a,d]*

5
6  [a] Department of Neurological Surgery, University of California, San Francisco, 675 Nelson
7  Rising Lane, San Francisco, California, 94158, USA
8  [b] Donders Institute for Brain, Cognition and Behaviour, Centre for Cognitive Neuroimaging,
9  Radboud University, Kapittelweg 29, Nijmegen, 6525 EN, The Netherlands
10 [c] Max Planck Institute for Psycholinguistics, Wundtlaan 1, Nijmegen, 6525 XD, Netherlands
11 [d] Weill Institute for Neurosciences, University of California, San Francisco, 675 Nelson
12 Rising Lane, San Francisco, California, 94158, USA
13
14 * Please address correspondence to: edward.chang@ucsf.edu
15
16
17

**ABSTRACT**

In speech, listeners extract continuously-varying spectrotemporal cues from the acoustic signal to perceive discrete phonetic categories. Spectral cues are spatially encoded in the amplitude of responses in phonetically-tuned neural populations in auditory cortex. It remains unknown whether similar neurophysiological mechanisms encode temporal cues like voice-onset time (VOT), which distinguishes sounds like /b/-/p/. We used direct brain recordings in humans to investigate the neural encoding of temporal speech cues with a VOT continuum from /ba/ to /pa/. We found that distinct neural populations respond preferentially to VOTs from one phonetic category, and are also sensitive to sub-phonetic VOT differences within a population's preferred category. In a simple neural network model, simulated populations tuned to detect either temporal gaps or coincidences between spectral cues captured encoding patterns observed in real neural data. These results demonstrate that a spatial/amplitude neural code underlies the cortical representation of both spectral and temporal speech cues.

**KEYWORDS**
speech perception; electrocorticography (ECoG); human auditory cortex; temporal processing; voice-onset time (VOT); categorical perception; sub-phonetic detail

## INTRODUCTION

During speech perception, listeners must extract acoustic cues from a continuous sensory signal and map them onto discrete phonetic categories, which are relevant for meaning[1, 2]. Many such cues to phonological identity are encoded within the fine temporal structure of speech[3–5]. For example, voice-onset time (VOT), defined as the interval between a stop consonant's release and the onset of vocal fold vibration (acoustically, the *burst* and the *voicing*), is a critical cue that listeners use to distinguish *voiced* (e.g., /*b*/, /*d*/, /*g*/) from *voiceless* (e.g., /*p*/, /*t*/, /*k*/) stop consonants in English[6, 7]. When the burst and voicing are roughly coincident (short VOT; ~0ms), listeners perceive a bilabial stop as a /*b*/, but when voicing follows the burst after a temporal gap (long VOT; ~50ms), listeners hear a /*p*/.

Recent evidence from human electrocorticography (ECoG) has shown that information about a speech sound's identity is encoded in the amplitude of neural activity at phonetically-tuned cortical sites in the superior temporal gyrus (STG)[8]. Distinct neural populations in this region respond selectively to different classes of phonemes that share certain spectral cues, such as the burst associated with stop consonants or the characteristic formant structure of vowels produced with specific vocal tract configurations. However, it is unclear whether phonetic categories distinguished by temporal cues (e.g., voiced vs. voiceless stops) are represented within an analogous spatial encoding scheme. If so, this would entail that local neural populations are tuned to detect not merely the presence of certain spectral cues (the burst and voicing), but also their timing relative to one another.

In addition to distinguishing phonetic categories, the exact VOT of a given utterance of a /*b*/ or a /*p*/ will vary considerably depending on numerous factors such as speech rate, phonetic context, and speaker accent[9–15]. Although only categorical phonetic identity (e.g., whether a particular VOT is more consistent with a /*b*/ or a /*p*/) is strictly necessary for understanding meaning, sensitivity to fine-grained sub-phonetic detail (e.g., whether a particular /*p*/ was pronounced with a 40ms vs. a 50ms VOT) is also crucial for robust speech perception, allowing listeners to flexibly adapt and to integrate multiple cues to phonetic identity online in noisy, unstable environments[16–21]. However, the neurophysiological mechanisms that support listeners' sensitivity[22–28] to such detailed speech representations are not known. We tested whether sub-phonetic information might be encoded in the neural response amplitude of the same acoustically-tuned neural populations that encode phonetic information in human auditory cortex.

To address these questions, we recorded neural activity directly from the cortex of seven human participants using high-density ECoG arrays while they listened to and categorized syllables along a VOT continuum from /*ba*/ (0ms VOT) to /*pa*/ (50ms VOT). We found that the amplitude of cortical responses in STG simultaneously encodes both phonetic and sub-phonetic information about a syllable's initial VOT. In particular, spatially discrete neural populations respond preferentially to VOTs from one category (either /*b*/ or /*p*/). Furthermore, peak response amplitude is modulated by stimulus VOT within each population's preferred – but not its non-preferred – voicing category (e.g., stronger response to 0ms than to 10ms VOT in voiced-selective [/*b*/-selective] neural populations). This same encoding scheme emerged in a computational neural network model simulating neuronal populations as leaky integrators tuned to detect either temporal coincidences or gaps between distinct spectral cues. Our results provide direct evidence that phonetic and sub-phonetic information carried by VOT are represented within spatially discrete, phonetically-tuned neural populations that integrate temporally-

84 distributed spectral cues in speech. This represents a crucial step towards a unified model of
85 cortical speech encoding, demonstrating that both spectral and temporal cues and both phonetic
86 and sub-phonetic information are represented by a common (spatial) neural code.
87
88 **RESULTS**
89
90     Participants listened to and categorized speech sounds from a digitally synthesized
91 continuum of consonant-vowel syllables that differed linearly only in their voice-onset time
92 (VOT) from /*ba*/ (0ms VOT) to /*pa*/ (50ms VOT). This six-step continuum was constructed by
93 manipulating only the relative timing of the spectral burst and the onset of voicing while holding
94 all other acoustic properties of the stimuli constant (**Figures 1A/B**; see **Methods**)(29). Analysis
95 of participants' identification behavior confirmed that stimuli with longer VOTs were more often
96 labeled as /*pa*/ (mixed effects logistic regression: $\beta_{VOT} = 0.19$, $t = 17.78$, $p = 5.6*10^{-63}$; data for
97 example participant in **Figure 1C**; data for all participants in **Figure 1-figure supplement 1**).
98 Moreover, and consistent with past work, listeners' perception of the linear VOT continuum was
99 sharply non-linear, a behavioral hallmark of categorical perception(30–32). A psychophysical
100 category boundary between 20ms and 30ms divided the continuum into stimuli most often
101 perceived as voiced (/*b*/: 0ms, 10ms, 20ms VOTs) or as voiceless (/*p*/: 30ms, 40ms, 50ms
102 VOTs).
103
104 **Temporal cues to voicing category are encoded in spatially distinct neural populations**
105
106     To investigate neural activity that differentiates the representation of speech sounds based
107 on a temporal cue like VOT, we recorded high-density electrocorticography in seven participants
108 while they listened to the VOT continuum. We examined high-gamma power (70-150 Hz)(33–
109 36), aligned to the acoustic onset of each trial (burst onset), at every speech-responsive electrode
110 on the lateral surface of the temporal lobe of each patient (n = 346 electrodes; see **Methods** for
111 details of data acquisition, preprocessing, and electrode selection).
112     We used nonparametric correlation analysis (Spearman's $\rho$) to identify electrodes where
113 the peak high-gamma amplitude was sensitive to stimulus VOT. Across all participants, we
114 found 49 VOT-sensitive sites, primarily located over the lateral mid-to-posterior STG,
115 bilaterally. Peak response amplitude at these VOT-sensitive electrodes reliably discriminated
116 between voicing categories, exhibiting stronger responses to either voiced (/*b*/; VOT = 0-20ms; n
117 = 33) or voiceless (/*p*/; VOT = 30-50ms; n = 16) stimuli (**Figure 1D**; locations of all sites shown
118 in **Figures 2A** and **1-figure supplement 2**). We observed that, within individual participants,
119 electrodes spaced only 4mm apart showed strong preferences for different voicing categories,
120 and we did not observe any clear overall regional or hemispheric patterns in the prevalence or
121 selectivity patterns of VOT-sensitive electrodes (see **Methods** for additional information).
122     Robust category selectivity in voiceless-selective (V-) and voiced-selective (V+) neural
123 populations emerged as early as 50-150ms post-stimulus onset and often lasted for several
124 hundred milliseconds (example electrodes in **Figure 1E**). Across all VOT-sensitive electrodes,
125 voicing category selectivity was reliable whether a trial's voicing category was defined based on
126 the psychophysically-determined category boundary (0-20ms vs. 30-50ms VOTs; V- electrodes:
127 $z = 3.52$, $p = 4.4x10^{-4}$; V+ electrodes: $z = -5.01$, $p = 5.4x10^{-7}$; Wilcoxon signed-rank tests) or
128 based on the actual behavioral response recorded for each trial (V- electrodes: $p = 4.9x10^{-4}$; V+
129 electrodes: $p = 6.1x10^{-5}$; Wilcoxon signed-rank tests).

130      These results show that spatially distinct neural populations in auditory cortex are tuned
131 to speech sound categories defined by a temporal cue. Critically, if individual neural populations
132 only responded to spectral features (e.g., to the burst or to the onset of voicing), we would not
133 have observed overall amplitude differences in their responses to /*b*/ versus /*p*/ categories.
134      Given this pattern of spatial tuning, we tested whether the voicing category of single
135 trials could be reliably decoded from population neural activity across electrodes. For each
136 participant, we trained a multivariate pattern classifier (linear discriminant analysis with leave-
137 one-out cross validation) to predict trial-by-trial voicing category using high-gamma activity
138 across all speech-responsive electrodes on the temporal lobe during the peak neural response
139 (150-250ms after stimulus onset; see **Methods**). We found that, across participants, classification
140 accuracy was significantly better than chance (Wilcoxon signed-rank test: $p = 0.016$; **Figure 1F**,
141 leftmost box plot), demonstrating that spatially and temporally distributed population neural
142 activity during the peak response contains information that allows for decoding of a temporally-
143 cued phonetic distinction in speech.

145 **Peak neural response amplitude robustly encodes voicing category**

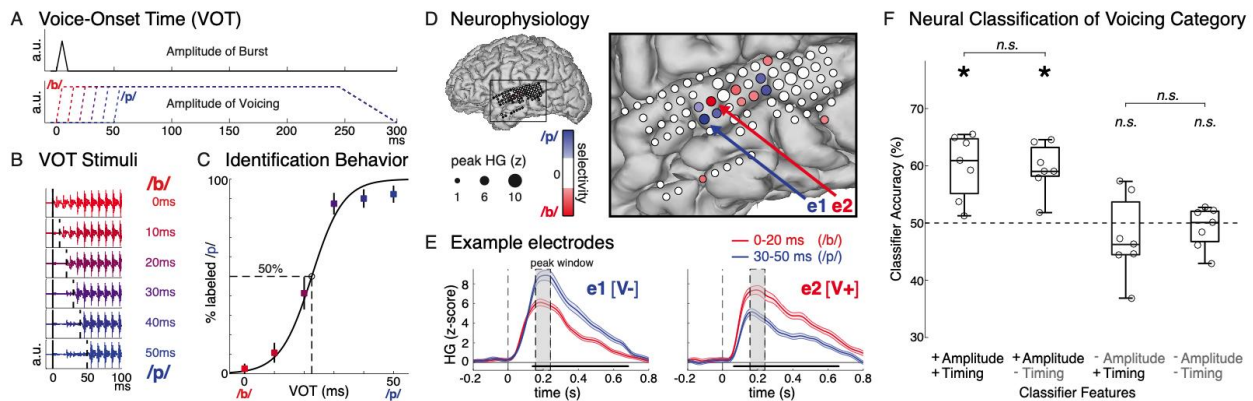147      Next, we asked which features of the population neural response encode voicing
148 category. Specifically, we evaluated three alternatives for how temporally-cued voicing category
149 is encoded by high-gamma responses in cortex during the peak neural response: (1) the spatial
150 pattern of peak response amplitude across electrodes, (2) the temporal patterns of evoked
151 responses across electrodes during the peak response, or (3) both amplitude and timing of neural
152 activity patterns. We tested these hypotheses by selectively corrupting amplitude and/or temporal
153 neural features that were inputs for the classifier. As with the previous analyses, and following
154 prior work on speech sound encoding(8), these analyses (**Figure 1F**) focused on cortical high-
155 gamma activity during the peak response window (150-250ms after stimulus onset; but see
156 **Figure 3** for analyses of an earlier time window).
157      To corrupt temporal information, we randomly jittered the exact timing of the neural
158 response for each trial by shifting the 100ms analysis window by up to ±50ms. Because the
159 uniform random jitter was applied independently to each trial, this procedure disrupts any
160 temporal patterns during the peak neural response that might reliably distinguish trials of
161 different voicing categories, such as precise (millisecond-resolution) timing of the peak response
162 at an electrode or the dynamics of the evoked response during the peak window, including *local*
163 temporal dynamics (during a single electrode's peak response) or *ensemble* temporal dynamics
164 (the relative timing of responses of spatially-distributed electrodes in the same participant). To
165 corrupt amplitude information, we eliminated any condition-related differences in the peak
166 response amplitude at every electrode. For each electrode, the evoked high-gamma response to
167 all trials within a given voicing category were renormalized so that the average responses to both
168 voicing categories had identical amplitudes at the peak, but could still vary reliably in the timing
169 and dynamics during the peak window. These techniques allowed us to examine the relative
170 contributions of temporal and amplitude information contained within the peak neural response
171 window to the classification of voicing category (see **Methods** for detailed description of this
172 approach).
173      Across participants, we found that, when the classifiers had access to amplitude
174 information but not timing information (+Amplitude/-Timing) during the peak response,
175 performance was significantly better than chance (Wilcoxon signed-rank test: $p = 0.016$; **Figure**

**1F**). Furthermore, despite the profound corruption of temporal information in the neural responses, classification accuracy was statistically comparable to the model that had access to both amplitude and timing information (+Amplitude/+Timing; Wilcoxon signed-rank test: $p = 0.69$; **Figure 1F**), suggesting that amplitude information alone is sufficient for classifying a trial's voicing category.

In contrast, when amplitude information was corrupted and only temporal patterns in the peak response window were reliable (-Amplitude/+Timing), classifier performance was not different from chance (Wilcoxon signed-rank test: $p = 0.69$; **Figure 1F**) and was worse for every participant compared to the model with both types of information (Wilcoxon signed-rank test: $p = 0.016$). Finally, we compared the model with only timing information to a model where both amplitude and timing information during the peak window were corrupted (-Amplitude/-Timing). We found that preserving timing information alone had no effect on classification performance compared to the most impoverished model (-Amplitude/-Timing; Wilcoxon signed-rank test: $p = 0.58$; **Figure 1F**), which also failed to perform better than chance (Wilcoxon signed-rank test: $p = 0.94$; **Figure 1F**). Together, these results constitute evidence for a spatial/amplitude code for speech categories that differ in a temporal cue. Thus, localized peak high-gamma response amplitude spatially encodes voicing of single trials in STG, analogous to other spectrally-cued phonetic features(8). Note that, while spatial (and not temporal) patterns of high-gamma responses robustly encode voicing during this critical peak window, we later describe additional analyses that address possible temporal encoding patterns in the local field potential (**Figure 1-figure supplements 3** and **4**) and in an earlier time window (**Figure 3**).



*Fig. 1. Speech sound categories that are distinguished by a temporal cue are spatially encoded in the peak amplitude of neural activity in distinct neural populations. A. Stimuli varied only in voice-onset time (VOT), the duration between the onset of the burst (**top**) and the onset of voicing (**bottom**) (a.u. = arbitrary units). B. Acoustic waveforms of the first 100ms of the six synthesized stimuli. C. Behavior for one example participant (mean ± bootstrap SE). Best-fit psychometric curve (mixed effects logistic regression) yields voicing category boundary between 20-30ms (50% crossover point). D. Neural responses in the same representative participant show selectivity for either voiceless or voiced VOTs at different electrodes. Electrode size indicates peak high-gamma (HG; z-scored) amplitude at all speech-responsive temporal lobe sites. Electrode color reflects strength and direction of selectivity (Spearman's ρ between peak HG amplitude and VOT) at VOT-sensitive sites (p < 0.05). E. Average HG responses (± SE) to voiced (0-20ms VOTs; red) and voiceless (30-50ms VOTs; blue) stimuli in two example electrodes from D, aligned to stimulus onset (e1: voiceless-selective, V-; e2: voiced-selective, V+). Horizontal black bars indicate timepoints with category discriminability (p < 0.005). Grey boxes mark average peak window (± SD) across all VOT-sensitive electrodes (n = 49). F. Population-based classification of*

*voicing category (/p/ vs. /b/) during peak window (150-250ms after stimulus onset). Chance is 50%. Boxes show interquartile range across all participants; whiskers extend to best- and worst-performing participants; horizontal bars show median performance. Asterisks indicate significantly better-than-chance classification across participants (p < 0.05; n.s. = not significant). Circles represent individual participants.*

The encoding of stop consonant voicing in the amplitude of evoked high-gamma responses in STG suggests that the representation of temporally-cued phonetic features may be explained within the same neural coding framework as the representation of spectrally-cued phonetic features. However, previous work on the cortical representation of voicing has identified a role for temporal information in the local field potential (LFP) (37, 38), which is dominated by lower- frequencies (39, 40).

To link our results with this existing literature, we conducted a series of exploratory analyses of the neural responses to our stimuli using the raw voltage (LFP) signal. For each VOT-sensitive electrode (defined in the high-gamma analysis), we estimated the correlations between VOT and peak latency and between VOT and peak amplitude for 3 peaks in the auditory evoked potential (AEP) occurring approximately 75-100 ms ($P_\alpha$), 100-150 ms ($N_\alpha$), and 150-250 ms ($P_\beta$) after stimulus onset (**Figure 1-figure supplement 3**)(41, 42). We found that some VOT-sensitive electrodes encoded VOT in the latency of these peaks (e.g., **Figure 1-figure supplement 4**, **panels E/I/M**), replicating previous results (43). However, among electrodes that encode VOT in peak high-gamma amplitude, there exist many more electrodes that *do not* encode VOT in these temporal features of the AEP, and many that also encode VOT in the amplitude of these AEP peaks (**Figure 1-figure supplements 3** and **4**). This further supports the prominent role that amplitude information plays in the neural representation of voicing and VOT, both in high-gamma and in the LFP. Therefore, subsequent analyses focus on the high-gamma amplitude. (For detailed descriptions of these LFP analyses and their results, see **Methods** and **Figure 1-figure supplements 3** and **4**.)

**Peak response amplitude encodes sub-phonetic VOT information within preferred category**

Next, we assessed whether VOT-sensitive neural populations (**Figure 2A**), which reliably discriminate between phonetic categories (voiced vs. voiceless), also encoded within-category sub-phonetic detail in the peak response amplitude. Specifically, the cortical representation of stimuli from the same voicing category but with different VOTs (e.g., 30, 40, and 50ms VOTs that all correspond to /p/) could be either categorical (i.e., all elicit the same peak response amplitude) or graded (i.e., peak response amplitude depends on within-category VOT).
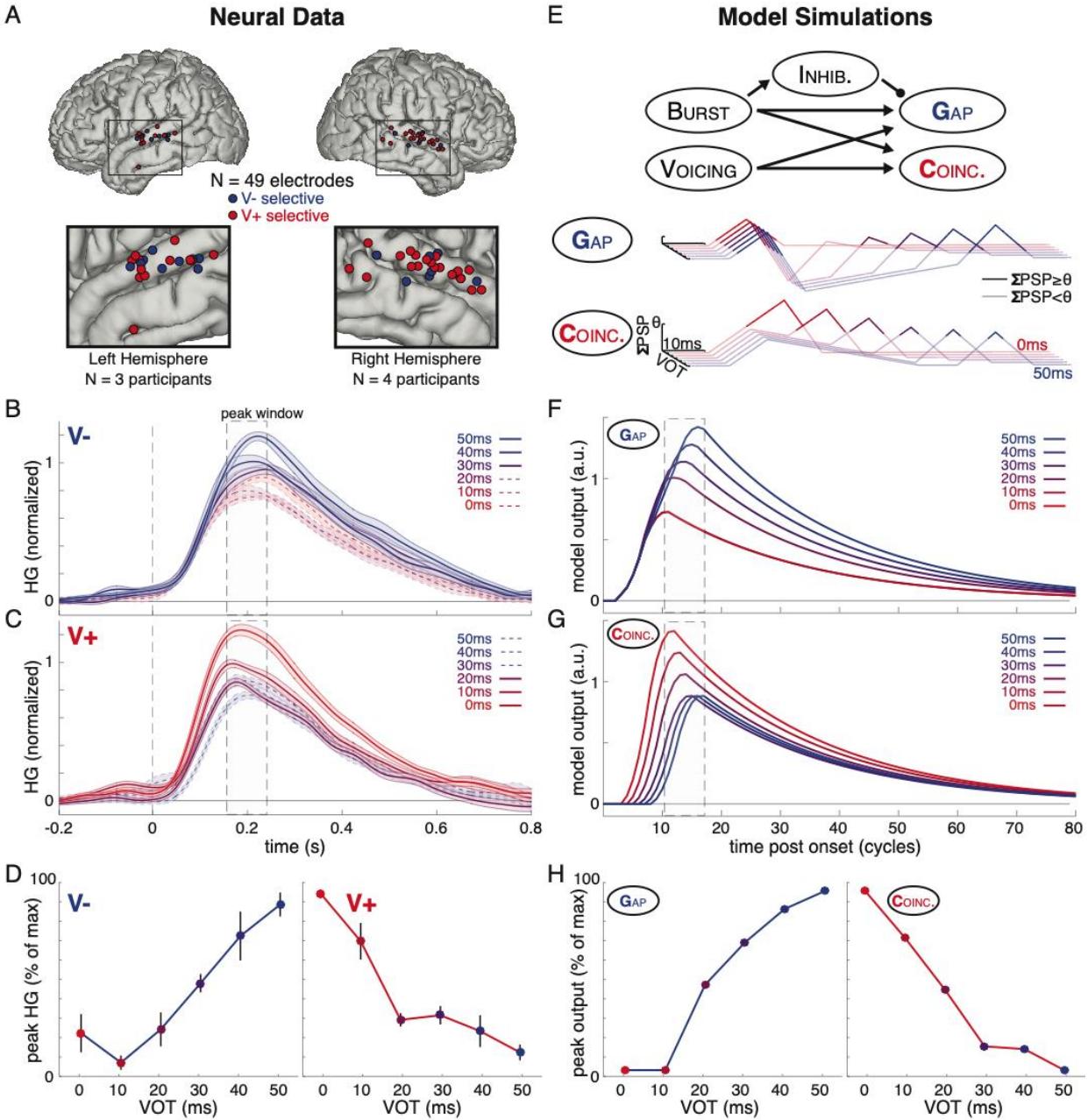
We examined the average responses to each of the six VOTs separately in the voiceless-selective electrodes (V-; **Figure 2B**) and the voiced-selective electrodes (V+; **Figure 2C**). We observed clear differences in activity evoked by different VOTs at the peak response (~200ms after stimulus onset), even within the same voicing category, consistent with sensitivity to sub-phonetic detail(44–47). However, the discriminability of responses to within-category VOTs depended on the preferred voicing category of a given electrode.

To quantify this observation, at each electrode, we computed the rank-based correlation (Spearman's $\rho$) between stimulus VOT and peak response amplitude separately for each voicing category (0-20ms and 30-50ms VOTs). This procedure resulted in two correlation coefficients for each VOT-sensitive site ($\rho_{0-20}$, $\rho_{30-50}$) and corresponding test statistics reflecting the strength

of within-category amplitude encoding of stimulus VOT in each voicing category. These test statistics (one per voicing category per VOT-sensitive electrode) then served as the input data for a series of signed-rank statistical tests to assess overall within-category encoding properties of groups of electrodes (e.g., of all V- electrodes) (see **Methods** for details). For example, consider V- electrodes, which exhibit stronger responses, overall, for voiceless stimuli (30-50ms VOTs) compared to voiced stimuli (0-20ms VOTs). Across V- electrodes, we found that voiceless stimuli with longer VOTs (i.e., closer to the preferred category's 50ms endpoint VOT) also elicit increasingly stronger responses (Wilcoxon signed-rank test: $z = 3.52$, $p = 4.4\text{x}10^{-4}$). At the same V- sites, however, within-category VOT does not reliably predict response amplitude among (non-preferred) voiced stimuli (Wilcoxon signed-rank test: $z = -1.60$, $p = 0.11$; **Figure 2B**: differences among solid blue lines but not dashed red lines). Across all V- and V+ electrodes, peak high-gamma response amplitude encoded stimulus VOT within the preferred category (Wilcoxon signed-rank test: $z = 6.02$, $p = 1.7\text{x}10^{-9}$), but not the nonpreferred category (Wilcoxon signed-rank test: $z = 1.31$, $p = 0.19$). While V- electrodes encoded sub-phonetic VOT more robustly within the voiceless category than within the voiced category (**Figure 2D**, **left**; Wilcoxon signed-rank test: $z = 3.00$, $p = 2.7\text{x}10^{-3}$), the opposite pattern emerged for V+ electrodes, which encoded sub-phonetic VOT more robustly within the voiced category than within the voiceless category (**Figure 2D**, **right**; Wilcoxon signed-rank test: $z = 3.78$, $p = 1.6\text{x}10^{-4}$).

Together, these analyses revealed two key results: (1) VOT encoding in human STG is not purely categorical, but also (2) the relationship between response amplitude and VOT is not linear across the entire continuum (**Figure 2D**). These results suggest that, even at the level of STG, the brain maintains information about the specific, sub-phonetic details of individual speech sounds. The asymmetrical pattern of within-category encoding suggests that individual neural populations in human auditory cortex encode information about both the category identity of a speech sound and its more fine-grained acoustic properties, or its "category goodness."(22, 44, 48)

288
289 ***Fig. 2. Human auditory cortex encodes both phonetic (between-category) and sub-phonetic (within-***
290 ***category) information in peak response amplitude, which can be modeled by a simple neural network***
291 ***that implements temporal gap and coincidence detection.*** *A. Spatial distribution of VOT-sensitive*
292 *electrodes across all (on standardized brain).* ***B.*** *Average (±SE) normalized HG response to each VOT*
293 *across all voiceless-selective (V-) electrodes, aligned to stimulus onset. Line style denotes category*
294 *membership of a given VOT (solid: preferred category; dashed: non-preferred category). Grey box marks*
295 *average peak window (±SD) across all VOT-sensitive electrodes.* ***C.*** *Average (±SE) normalized response*
296 *to each VOT across all voiced-selective (V+) electrodes.* ***D.*** *Average (±SE) peak response to each VOT*
297 *stimulus for V- electrodes (****left****) and V+ electrodes (****right****) (see* ***Methods****).* ***E.*** *A simple neural network*
298 *model (****top****) comprised of five leaky integrator nodes was implemented to examine computational*
299 *mechanisms that could account for the spatial encoding of a temporal cue (VOT). Arrows and circle*
300 *represent excitatory and inhibitory connections between nodes. See* ***Methods*** *for details on model*

9

*parameters. Postsynaptic potentials (PSPs) illustrate the internal dynamics of the gap detector (**GAP**, middle) and coincidence detector (**COINC., bottom**) in response to simulated VOT stimuli (line color). Outputs (**panels F/G**) are triggered by suprathreshold instantaneous PSPs (ΣPSP≥θ, dark lines) but not by subthreshold PSPs (ΣPSP<θ; semitransparent lines). **F.** Model outputs (a.u. = arbitrary units) evoked by simulated VOT stimuli for **GAP** (1 cycle = 10ms). Note that outputs for 0ms and 10ms VOTs are overlapping. No error bars shown because model simulations are deterministic. Grey box marks average peak window (across **panels F/G**); width matches peak window of real neural data (**panels B/C**). **G.** Model outputs for **COINC.** **H.** Peak response to each simulated VOT stimulus for **GAP** (left) and **COINC.** (right).*

## A simple neural network model of VOT encoding in STG

Thus far, we have demonstrated that a temporal cue that distinguishes speech sounds is represented by a spatial/amplitude code(49, 50) in human STG. To understand how this could be implemented computationally in the brain, we built an architecturally minimalistic neural network (**Figure 2E**, **top**). The network was designed to implement a small set of basic computations, motivated by well-established models of temporal processing(51–57). Specifically, our model employs discrete integrator units that detect temporal gaps or coincidences between distinct spectral events by incorporating canonical neurophysiological mechanisms that allow current input to modulate a unit's sensitivity to subsequent input in highly specific ways.

The entire model is comprised of just five localist units: a burst detector, a voicing detector, a gap detector (**GAP**), a coincidence detector (**COINC.**), and an inhibitory unit. Conventional leaky integrator dynamics governed continuously varying activation values of each rectified linear unit within the model(58, 59), with the activity $a_i(t)$ of a given unit $i$ at time $t$ depending on its prior activity $a_i(t-1)$, the weighted sum of its excitatory and inhibitory inputs $\sum_j w_{ji} * a_j(t-1)$, and unit-specific activation parameters (e.g., propagation threshold [$\theta$], decay rate). To illustrate intuitively how time-dependent neuronal properties can give rise to spatially-localized temporal cue processing, model parameters and connection weights were set manually (see **Methods** for details; **Figure 2-figure supplement 1; Supplementary File 2**). We presented the network with simplified inputs mimicking the spectral and temporal properties of the six VOT stimuli used in the ECoG experiment (**Figure 1A**; see **Methods**; **Supplementary File 3**). Presentation of burst and voicing inputs triggered propagation of activation that spread through the network, and our analyses assessed how the resulting activation dynamics differed depending on VOT.

The simulated responses of **GAP** and **COINC.** to VOTs of 0-50ms are shown in **Figures 2F/G**. We observed striking qualitative similarities between **GAP**'s simulated outputs (**Figure 2F**) and the real neural responses of V- electrodes (**Figure 2B**), and between **COINC.**'s outputs (**Figure 2G**) and the V+ electrodes (**Figure 2C**). By design, voicing category is clearly distinguished in both **GAP** and **COINC.**, with **GAP** responding more strongly to longer (voiceless) VOTs (30-50ms), and **COINC.** responding more strongly to shorter (voiced) VOTs (0-20ms). This demonstrates that spatial encoding of temporal cues (gaps vs. coincidences) can arise naturally within a simple, biologically-inspired neural network(51–57).

Perhaps more surprisingly, we also found that both **GAP** and **COINC.** detector units exhibit sensitivity to within-category VOT distinctions (**Figure 2H**). These partially graded activations mirror the pattern observed in the neural data (**Figure 2D**), where V- electrodes and **GAP** units

10

347 are only sensitive to differences among long (voiceless) VOTs, and V+ electrodes and *COINC.*
348 units are only sensitive to differences among short (voiced) VOTs.

349 These relatively sophisticated dynamics are the natural result of well-established
350 computational and physiological mechanisms. Within the model, the burst and voicing detector
351 units are tuned to respond independently to distinct spectral cues in the simulated acoustic input.
352 Hence, the relative timing of their responses, but not their amplitudes, differ as a function of
353 VOT. Both the gap (*GAP*) and the coincidence (*COINC.*) detector units receive excitatory input
354 from both the burst and voicing detector units, but *GAP* and *COINC.* differ in how they integrate
355 these inputs over time. Specifically, as described below, while initial excitatory input (from the
356 burst detector) temporarily *decreases* the sensitivity of *GAP* to immediate subsequent excitatory
357 input (from the voicing detector), the opposite is true of *COINC.*

358 In particular, prior work has shown that one computational implementation of gap
359 detection involves configuration of a *slow inhibitory postsynaptic potential* (*IPSP*) microcircuit
360 (**Figure 2E**, **middle**)(51, 52, 60, 61). In our model, activity in the burst detector following burst
361 onset elicits fast suprathreshold *excitatory postsynaptic potentials* (*EPSPs*) in both *GAP* and the
362 inhibitory unit, immediately followed by a longer-latency ("slow") IPSP in *GAP*. This slow IPSP
363 renders *GAP* temporarily insensitive to subsequent excitatory input from the voicing detector,
364 meaning that voicing-induced excitation that arrives too soon (e.g., 10ms) after the burst input,
365 when inhibition is strongest, is not able to elicit a second suprathreshold EPSP in *GAP*.
366 Consequently, all short VOTs (below some threshold) elicit uniformly weak responses in *GAP*
367 that reflect only the initial excitatory response to the burst (see, e.g., indistinguishable responses
368 to 0ms and 10ms VOTs in **Figure 2F**). However, as *GAP* gradually recovers from the burst-
369 induced slow IPSP, later-arriving voicing input (i.e., longer VOTs) tends to elicit suprathreshold
370 responses that grow increasingly stronger with longer gaps, until *GAP* has reached its pre-IPSP
371 (resting) baseline. In this way, our implementation of gap detection naturally captures three key
372 patterns observed across V- electrodes (**Figure 2H**, **left**; **Figure 2D**, **left**): (1) amplitude
373 encoding of a temporally cued category (selectivity for gaps over coincidences); (2) amplitude
374 encoding of within-category differences in the preferred category (amplitude differences among
375 gaps of different durations); and (3) no amplitude encoding of differences within the non-
376 preferred category (uniformly lower amplitude responses to short VOTs of any duration).

377 In contrast, coincidence detection(54–56, 62–64) (**Figure 2E**, **bottom**) emerges in the
378 model because activity in the burst detector evokes only a subthreshold EPSP in *COINC.*,
379 temporarily increasing *COINC.*'s sensitivity to immediate subsequent excitatory input (from the
380 voicing detector). During this period of heightened sensitivity, voicing-induced excitatory input
381 that arrives simultaneously or after short lags can elicit larger amplitude (additive) EPSPs than
382 could voicing-induced excitatory input alone. Because the magnitude of the initial burst-induced
383 EPSP gradually wanes, the summation of EPSPs (from the burst and voicing) is greatest (and
384 hence elicits the strongest response) for coincident burst and voicing (0ms VOT), and the
385 magnitude of *COINC.*'s response to other voiced stimuli (e.g., 10-20ms VOTs) becomes weaker
386 as the lag between burst and voicing increases. Finally, in voiceless stimuli, since voicing arrives
387 late enough after the burst (30+ ms) that there is no residual boost to *COINC.*'s baseline post-
388 synaptic potential, elicited responses are entirely driven by a suprathreshold voicing-induced
389 EPSP that reaches the same peak amplitude for all voiceless stimuli. Thus, our implementation of
390 coincidence detection captures three key patterns observed in V+ electrodes (**Figure 2H**, **right**;
391 **Figure 2D**, **right**): (1) amplitude encoding of a temporally cued category (selectivity for
392 coincidences over gaps); (2) amplitude encoding of within-category differences in the preferred

category (amplitude differences among stimuli with short VOTs); and (3) no amplitude encoding of differences within the non-preferred category (uniformly lower amplitude responses to long VOTs of any duration).

In summary, the neurophysiological dynamics underlying local STG encoding of VOT can be modeled using a simple, biologically-inspired neural network. The computational model captures both the between-category (phonetic) and within-category (sub-phonetic) properties of observed neural representations via well-established physiological mechanisms for gap and coincidence detection(51–57).

**Mechanisms that explain local category-selectivity also predict early temporal dynamics**

Thus far, we have focused on the encoding of speech sounds that differ in VOT based on activity patterns around the peak of the evoked response. However, in comparing the real and simulated neural data (**Figure 2**), we also observed a qualitative resemblance with respect to the onset latencies of evoked responses. Specifically, the timing of the evoked neural responses (relative to burst onset) appeared to depend on stimulus VOT in V+ electrodes and in the coincidence detector (*COINC.*) unit (**Figures 2C/G**), but not in V- electrodes or in the gap detector (*GAP*) unit (**Figure 2B/F**). This pattern could suggest that early temporal dynamics of the evoked response contribute to the pattern of category selectivity observed at the peak.
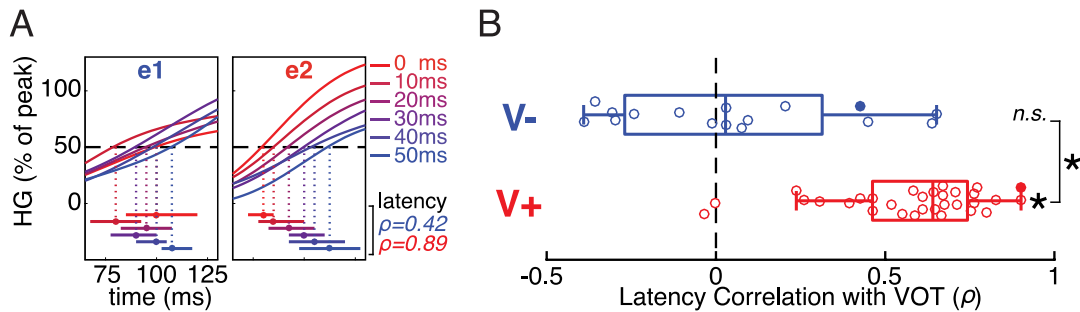
We examined the neural activity evoked by each VOT stimulus in V- and V+ electrodes at the onset of the response, typically beginning approximately 75-125ms after stimulus (burst) onset. In the same two example electrodes from **Figure 1E**, we observed clear differences in the relationship between response onset latency and VOT (**Figure 3A**). To quantify the onset latency for each electrode to each VOT stimulus, we found the first timepoint after stimulus onset where the evoked high gamma response exceeded 50% of the electrode's overall peak amplitude (grand mean across conditions). The rank correlation between VOT and response onset latency for e1 (a V- electrode) was substantially lower (Spearman's $\rho = 0.42$) than for e2 (a V+ electrode; $\rho = 0.89$).

A bootstrapped rank-based correlation coefficient was computed for each V- and V+ electrode (1000 resamples; see **Methods**). We found that response onset latency was strongly associated with VOT for V+, but not V-, electrodes (Wilcoxon signed-rank tests: V+, $p = 1.6 \times 10^{-6}$; V-, $p = 0.57$), and this difference between the two electrode types was highly reliable (Mann-Whitney rank-sum test: $p = 1.7 \times 10^{-5}$) (**Figure 3B**).

The association between VOT and response latency also differed in *GAP* versus *COINC.* units in the model simulations (**Figures 2F/G**), with VOT-dependent response latencies emerging for *COINC.*, but not *GAP*. Closer examination of the model's internal dynamics reveals how the same time-dependent mechanisms that give rise to peak amplitude encoding of VOT are also responsible for these early temporal dynamics. As described above, the category selectivity of *GAP* (voiceless) and *COINC.* (voiced) results from how each unit's subsequent activity is modulated after detection of the burst. While the burst always elicits a fast suprathreshold response in *GAP* (irrespective of VOT), *COINC.*'s response to the burst alone is subthreshold (**Figure 2E**, **middle** vs. **bottom**). Consequently, *GAP*'s initial response is evoked by the burst of any VOT stimulus, so the response onset latency (when aligned to burst onset) does not depend on VOT (**Figure 2F**). Conversely, *COINC.*'s earliest suprathreshold response is triggered by the onset of voicing, so the response onset latency (relative to burst onset) is later for longer VOTs (**Figure 2G**). Thus, the same well-established physiological mechanisms that give rise to peak

amplitude encoding of temporally-cued voicing categories also predict the early temporal dynamics we observe in real neural data.

Finally, **Figure 3** shows that, unlike during the peak response window (150-250ms after stimulus onset; **Figure 1F**), temporal information does encode VOT during an earlier window around the neural response onset in some neural populations. Indeed, both sub-phonetic and phonetic category-level information are carried by the onset latency of V+ electrodes, with evoked responses arising later at these sites for stimuli with progressively longer VOTs. Critically, the modeling results indicate that both the amplitude encoding patterns during the peak window and the temporal encoding patterns during the earlier onset window are captured by the same canonical neurophysiological mechanisms.



*Fig. 3. Early temporal dynamics of stimulus-evoked neural responses differ between voiceless-selective (V-) and voiced-selective (V+) electrodes. A. Normalized trial-averaged HG responses to each VOT stimulus (line color) in two example electrodes (e1 and e2; same electrodes shown in Figures 1D/E). The time window (x-axis) is relative to onset of the burst and precedes the peak response. Horizontal bars show estimates (bootstrapped mean ± SE) of response onset latency for each VOT (first timepoint exceeding 50% of electrode's average peak HG). Mean bootstrapped rank-based correlation (Spearman's ρ) between VOT and response onset latency shown for e1 (blue) and e2 (red). B. Across all V- electrodes, the bootstrapped correlation coefficients did not differ significantly from 0, suggesting that onset latency was time-locked to the burst. In contrast, across all V+ electrodes, the bootstrapped correlation coefficients were reliably positive (longer latencies for longer VOTs), and greater than for V- electrodes. Circles represent individual electrodes (filled: example electrodes in A). Boxes show interquartile range; whiskers extend to maximum/minimum of each group (excluding 2 outlier V+ electrodes); vertical bars are medians. Asterisks indicate significance ($p < 10^{-4}$; n.s. = not significant).*

**DISCUSSION**

This study investigated how voice-onset time (VOT), a temporal cue in speech, is represented in human auditory cortex. Using direct intracranial recordings, we found discrete neural populations located primarily on the bilateral posterior and middle STG that respond preferentially to either voiced sounds, where the onset of voicing is coincident with the burst or follows it after a short lag (20ms or less), or voiceless sounds, where the onset of voicing follows the burst after a temporal gap of at least 30-50ms.

Past work has also found that phonetic information about speech sounds is encoded in the amplitude of evoked neural responses at spatially localized cortical sites(8). In that work, however, STG activity was shown to encode the spectral properties of speech sounds most robustly, such as whether a phoneme is a vowel or a consonant and whether a consonant's spectrum is broadband (as in plosives, like /b/ and /p/) or is dominated by acoustic energy at high frequencies (as in fricatives, like /f/ and /s/).

The present results extend these earlier findings in a critical way, suggesting that the cortical representation of both spectral and temporal cues in speech follow a common spatial coding scheme. This result is also consistent with prior reports that neural response amplitude depends on VOT(8), but such results have often involved natural speech stimuli where voicing categories varied along many other spectral acoustic dimensions besides the temporal cue(65–68). Here, the digitally synthesized VOT stimuli were tightly controlled to vary only in the relative timing of two invariant spectral cues (burst and voicing), thereby demonstrating that this temporal speech cue is encoded in the peak high-gamma response amplitude of spatially distinct neural populations in human STG.

While the present results clearly implicate a spatial/amplitude code in the cortical representation of VOT, other work has described VOT-dependent temporal response patterns that can also be used to encode voicing categories(69–71). For instance, Steinschneider and colleagues have observed neurons and neuronal populations in primate and human auditory cortices in which short VOTs elicit a single-peaked neural response, while longer VOTs elicit a double-peaked response(37, 38, 43, 72–75). Under this "local" temporal coding model, the precise temporal dynamics of the response evoked at a single cortical site could distinguish voiced from voiceless VOTs. Our examination of the timing and amplitude of three peaks in the auditory evoked local field potentials of VOT-sensitive electrodes confirmed that such patterns do appear in some electrodes (**Figure 1-figure supplements 3** and **4**), clearly demonstrating that temporal and amplitude codes for VOT are not mutually exclusive (see also temporal encoding patterns in onset latencies of V+ electrodes; **Figure 3**). However, as with spectrally-defined phonetic contrasts (e.g., plosive vs. fricative(8)), it clear that the amplitude of the peak high-gamma (and, in many cases, of the LFP) response emerged as a robust representation of voicing category and of VOT.

VOT could also be encoded in the relative timing of responses in spatially-distributed, spectrally-tuned burst- and voicing-selective neural populations. Under this "ensemble" temporal coding model(76, 77), the pattern of neural activity evoked by voiced VOTs (characterized by roughly coincident burst and voicing cues) would differ from the pattern evoked by voiceless VOTs in the precise temporal latency of the response in a vowel-selective neural population (a voicing detector) compared to the response in a plosive-selective neural population (a burst detector). However, the fact that we found cortical sites in every participant that exhibited robust category-dependent differences in their peak response amplitude rules out the possibility that at

511 least these neural populations are merely responding to spectral cues in the burst or voicing
512 alone.
513     Notably, if either (or both) of these models – a local or ensemble temporal code – were
514 primarily responsible for the neural representation of VOT in the high-gamma range, then the
515 selective corruption of temporal information in a classifier (**Figure 1F**) should have reduced
516 neural decoding of voicing category to chance levels, while corrupting peak amplitude
517 information should have had little or no effect. We found the opposite pattern of results:
518 corrupting peak amplitude information had a devastating effect on the decoding of voicing
519 category, while corrupting the fine temporal patterns that could have discriminated between
520 voicing categories had no measurable impact on classifier performance. To be clear, our work
521 does not rule out the possibility that local or ensemble temporal codes may also play a role in the
522 cortical representation of VOT. However, it does highlight spatially-localized peak neural
523 response amplitude as a robust code for VOT. Thus, in contrast to prior work theorizing parallel,
524 but fundamentally different, coding schemes for spectrally- and temporally-cued phonetic
525 features(37, 38), we demonstrate evidence for a shared representation of both by high-gamma in
526 the human superior temporal lobe.
527     In order to explicitly test potential computational and physiological mechanisms that
528 could give rise to the observed spatial coding scheme, we implemented an architecturally simple
529 neural network model. Although it is well known that spectral information is represented by a
530 spatial neural code from the earliest stages of auditory transduction in the cochlea(78, 79), the
531 emergence of a spatial code for the representation of temporally-distributed cues in a transient
532 acoustic signal poses a nontrivial computational problem. Our model highlights one
533 parsimonious approach by which selectivity for either temporal gaps or coincidences could be
534 implemented by biologically-inspired neurophysiological microcircuits(51–57).
535     We found that, just like in the neural data, gap and coincidence detector units responded
536 to simulated voiced (/*b*/) and voiceless (/*p*/) stimuli with different response amplitudes. As such,
537 we need not invoke any specialized temporal code to account for the representation of temporally
538 cued phonetic features. Rather, our results provide evidence implicating a common neural coding
539 scheme in the neural representation of behaviorally relevant speech features, whether they are
540 embedded within the instantaneous spectrum or the fine temporal structure of the speech signal.
541 Recent ECoG evidence suggests an even more expansive view of the fundamental role of spatial
542 coding in cortical speech representation(80) in which different neural populations also encode
543 pitch(81) and key properties of the speech envelope such as onsets and auditory edges(82, 83).
544     Crucially, although the neural network was only designed to discriminate between
545 categories (i.e., gaps vs. coincidences), we also observed graded amplitude differences in
546 response to different VOTs (**Figure 2H**), but only in an electrode's preferred category. These
547 within-category patterns emerged naturally from the same computational properties that allowed
548 the network to capture basic between-category encoding: (1) the relative responsiveness of each
549 temporal integrator unit (*GAP*, *COINC.*) to its various inputs (burst, voicing, and inhibition); (2)
550 the time-dependent properties inherent to neuronal activation dynamics (e.g., decay of
551 postsynaptic potentials towards a unit's resting activation level); and (3) the nonlinear
552 transformation of postsynaptic inputs into response outputs (rectified linear activation function
553 controlled by a unit's propagation threshold).
554     This asymmetric within-category encoding scheme closely resembled the pattern
555 observed in real neurophysiological data, where peak response amplitude to VOTs within the
556 same voicing category only differed within a neural population's preferred category (**Figure**

557    **2D**). This result clearly demonstrates that human nonprimary auditory cortex maintains a robust,
558    graded representation of VOT that includes the sub-phonetic details about how a particular
559    speech token was pronounced(44–47). Even though sub-phonetic information is not strictly
560    necessary for mapping sound to meaning in stable, noise-free listening environments, this fine-
561    grained acoustic detail has demonstrable effects on listeners' behavior(22–28), and modern
562    theories of speech perception agree that perceptual learning (e.g., adaptation to accented
563    speakers) and robust cue integration would be impossible if the perception of speech sounds
564    were strictly categorical(16–20, 84–87). Crucially, these data suggest that the same
565    spatial/amplitude code that is implicated in the representation of *phonetic* information (from
566    spectral or temporal cues) can also accommodate the representation of *sub-phonetic* information
567    in the speech signal.
568         The onset latency results (**Figure 3**) established an entirely novel correspondence
569    between the real and simulated results that extended beyond the peak response window.
570    Response onset latencies of V- electrodes were time-locked to the burst (**Figures 2B** and **3**),
571    while responses of V+ electrodes were time-locked to voicing onset (**Figures 2C** and **3**). These
572    highly reliable neurophysiological results neatly match specific predictions of our parsimonious
573    model without the need to postulate additional mechanisms (**Figures 2F/G**).
574         The correspondence between simulated and real neural data in the onset latency results
575    may also have implications for the question of whether the observed temporal integration is
576    occurring locally in STG or is inherited from earlier levels of auditory processing (e.g., from
577    midbrain or primary auditory cortex). The model's gap and coincidence detectors (*GAP*, *COINC.*)
578    are designed to directly simulate neural populations in the STG. Their inputs from the burst and
579    voicing detectors are only spectrally processed, so, in the model, the temporal onset latency
580    dynamics (**Figures 2F/G**) first arise in *GAP* and *COINC.* As such, the fact that the model's
581    prediction is borne out in the neural data in STG (**Figures 2B/C** and **3**) is consistent with local
582    temporal integration in STG. While these modeling results do not definitively rule out temporal
583    integration at lower levels of the ascending auditory pathway, its potentially local emergence in
584    high-order auditory cortex illustrates how even relatively simple computational models can be
585    used to generate novel hypotheses, which can ultimately be tested in real neurophysiological
586    data.
587         Overall, the results of these model simulations illustrate how the same network properties
588    that transform temporal cues into a spatial code are also able to naturally explain at least three
589    additional patterns observed within category-selective neural populations: (1) the graded
590    encoding of VOT within a population's preferred category; (2) the lack of graded encoding of
591    VOT within a population's non-preferred category; and (3) the early temporal dynamics of
592    neural responses, which depend on a population's category-selectivity. Thus, the model provides
593    an explicit, mathematical account of multiple seemingly disparate observations about the
594    neurophysiological data, all of which arise directly from a parsimonious implementation of gap-
595    and coincidence-detection with well-established, theoretically-motivated neuronal circuits.
596         The model we present is just one of many possible architectures that could capture these
597    interesting properties of the neural response. For example, mechanisms like temporal delay lines
598    (54, 56) could also be used to implement gap detection. Broadly, we chose to implement a
599    simple hand-tuned neural network model to maximize our ability to explore the detailed
600    dynamics we observed in the neural data. Our approach follows a rich history of using these
601    types of hand-tuned models to explain a wide array of cognitive and perceptual phenomena
602    (including the perception of VOT in speech), as exemplified by the influential TRACE model of

603  speech perception(84). An alternative approach to modeling VOT perception is to train a neural
604  network to distinguish voiced from voiceless sounds based on distributed activation dynamics
605  within biologically-grounded spectral processing maps(88). Our model borrows aspects of these
606  two approaches (hand-tuning; biological plausibility) and it extends this past work by directly
607  modeling the time-dependent mechanisms that could give rise to continuously-varying neural
608  responses in STG.
609        While the model captured several notable features of the neural data (including some for
610  which it was not explicitly designed), we observed at least one inconsistency between the
611  simulated and real neural responses. The model predicted VOT-dependence in the latency of the
612  *peak* response in both *GAP* and *COINC.* units (**Figures 2F/G**), but we did not find evidence for
613  these fine-grained patterns in the high-gamma data (**Figures 2B/C**; see also lack of category-
614  dependent temporal patterns during peak window: **Figure 1F**). However, it is unclear whether
615  this is a false prediction of the model, or whether we did not observe the effect in the neural data
616  because of, for example, poor signal-to-noise ratio for this effect. Regardless of whether the
617  discrepancy arises from the model or the real data, it represents a gap in our mechanistic
618  understanding of the processing of this phenomenon, and should therefore be a target for further
619  research.
620        Although topographic functional organization is pervasive among many spatial neural
621  coding schemes described in sensory neuroscience, including for the representation of spectral
622  and temporal acoustic cues in audition (e.g., tonotopy in mammalian auditory cortex(78, 79) or
623  chronotopy in bats(89, 90)), this functional organization seems not to extend to the spatial code
624  for speech on the lateral temporal cortex in humans. As with tuning for spectrally-cued phonetic
625  features(8, 82) (e.g., plosives vs. fricatives), VOT-sensitive neural populations in the present
626  study were scattered throughout posterior and middle superior temporal gyrus with no
627  discernible topographical map of selectivity or evidence for lateralized asymmetries(71, 91),
628  although data limitations prevent us from ruling out this possibility entirely (for detailed results,
629  see **Methods**).
630        Most of the present analyses focused on the high-gamma component of the neural
631  response, but this work does not discount a potential role for lower-frequency oscillations in
632  speech perception(92, 93) or in the perception of phonemes(94, 95). Indeed, it is clear from the
633  exploratory analyses of auditory evoked local field potentials (**Figure 1-figure supplements 3**
634  and **4**) that there do exist complex associations between VOT and the amplitude/temporal
635  information carried in lower-frequency ranges. Future work should systematically investigate the
636  relationship between high-gamma and other neural signals (such as the local field potential),
637  their relative contributions to the perceptual experience of and neural representation of speech,
638  and the importance of detailed temporal information in each (see, e.g., 42).
639        Finally, it is critical to distinguish our results from studies describing neural correlates of
640  <u>categorical</u> speech perception, *per se* (e.g., 96). Neural responses to different VOT tokens that
641  are members of the same voicing category can only be considered truly categorical if the
642  responses are indiscriminable (e.g., 30, 97). In our results, acoustically distinct members of the
643  same phonetic category <u>*are*</u> distinguishable in neural populations that are selective for that
644  voicing category (**Figure 2**). In light of this graded VOT representation, the present results are
645  best interpreted as elucidating neural mechanisms of category perception, but not necessarily
646  categorical perception, of voiced vs. voiceless stop consonants. While limited coverage beyond
647  the superior temporal lobe precludes us from ruling out the influence of top-down categorical
648  perception (98–100) (possibly originating in frontal cortex (101–104)) on our results, it is notable

649 that the model we present (which does not posit top-down effects) suggests that top-down effects
650 may not be a necessary condition for explaining the observed non-linear encoding patterns (see
651 also 84, 85, 105–107).
652       In conclusion, the present results show that spatially-discrete neural populations in human
653 auditory cortex are tuned to detect either gaps or coincidences between spectral cues, and these
654 sites simultaneously represent both phonetic and sub-phonetic information carried by VOT, a
655 temporal speech cue found in almost all languages(7, 108). This demonstrates a common
656 (spatial) neural code in STG that accounts for the representation of behaviorally relevant
657 phonetic features embedded within the spectral and temporal structure of speech. From a simple
658 model that transforms a temporal cue into a spatial code, we observed complex dynamics that
659 show how a highly variable, continuous sensory signal can give rise to partially abstract, discrete
660 representations. In this way, our findings also add to a growing body of work highlighting the
661 critical role of human STG as a sensory-perceptual computational hub in the human speech
662 perception system(80, 81, 96, 102, 109–112).
663

**METHODS**

**Data and code availability.** All data and code associated with this study and necessary for replication of its results are available under a Creative Commons license at the associated Open Science Framework project page (https://osf.io/9y7uh/).(113)

**Participants.** A total of seven human participants with self-reported normal hearing were implanted with high-density (128 or 256 electrodes; 4 mm pitch) multi-electrode cortical ECoG surface arrays as part of their clinical treatment for epilepsy. Placement of electrode arrays was determined based strictly on clinical criteria. For all patients who participated in this study, coverage included peri-Sylvian regions of the lateral left (n = 3) or right (n = 4) hemisphere, including the superior temporal gyrus (STG). All participants gave their written informed consent before the surgery and affirmed it at the start of each recording session. The study protocol was approved by the University of California San Francisco Committee on Human Research. Data from two additional participants were excluded from analyses because of excessive epileptiform activity (artifacts) during recording sessions.

**Imaging.** Electrode positions (**Figure 1D** and **Figure 1-figure supplement 2**) were determined from post-surgical computed tomography (CT) scans and manually co-registered with the patient's MRI. Details of electrode localization and warping to a standardized brain (MNI; **Figure 2A**) are described elsewhere(114).

**Stimuli.** Stimuli (**Figure 1B**) were generated with a parallel/cascade Klatt-synthesizer KLSYN88a using a 20-kHz sampling frequency (5ms frame width in parameter tracks). All stimulus parameters were identical across stimuli, with the exception of the time at which the amplitude of voicing began to increase (in 10ms steps from 0ms to 50ms after burst onset; **Figure 1A**). The total duration of each stimulus was 300ms regardless of VOT. The onset noise-burst was 2ms in duration and had constant spectral properties across all stimuli. The dominant frequency ranges for the vowel were: F0 = 100 Hz; F1 = 736 Hz; F2 = 1221 Hz; F3 = 3241 Hz (consistent with a vocal tract length of 13.5 cm). Formant transitions always began at 30ms. The vowel's amplitude began ramping down 250ms after stimulus onset. The stimuli are made available among this study's supplementary materials and at the associated Open Science Framework page.(113)

**Behavioral Procedure.** During ECoG recording, the VOT stimuli were presented monaurally over free-field loudspeakers at a comfortably listening level via a custom MATLAB script(113) in a blocked pseudorandom order. Four of seven participants simultaneously performed a behavioral task wherein they indicated on each trial whether they heard "ba" or "pa" using a touchscreen tablet (programmed using a custom MATLAB GUI). In these recording sessions, the onset of the next trial began 500ms after a response was registered or 5 seconds after the end of the stimulus (if no response was registered). In sessions where participants chose to listen to the stimuli passively (instead of participating in the behavioral task), the onset of the next trial began approximately 1000ms after the end of the previous trial. **Supplementary File 1** reports number of trials per participant.

709 **Behavioral Analysis.** For the four participants who participated in the behavioral identification
710 task, individual trials were excluded from behavioral analysis if a participant did not make a
711 response or if the participant's reaction time was more than 3 standard deviations from the
712 participant's mean reaction time.
713
714 Behavioral response data were submitted to mixed effects logistic regression with a fixed effect
715 of VOT (coded as a continuous variable) and random intercepts for participants, allowing
716 individual participants to vary in their voicing category boundary. Using the best-fit model
717 estimates, we calculated the overall voicing category boundary across all participants ($\chi =$
718 21.0ms; **Figure 1-figure supplement 1**, **panel A**) and in the each individual participant (after
719 adjusting for random intercept fit for each participant; **Figure 1-figure supplement 1**, **panel B**,
720 and **Figure 1C**) as follows(115), where $\beta_0$ is the best-fit intercept and $\beta_{VOT}$ is the best-fit effect
721 of slope:

$$\chi = -\frac{\beta_0}{\beta_{VOT}}$$

722
723 **ECoG signal processing.**
724      **Recording and preprocessing.** Voltage fluctuations were recorded and amplified with a
725 multichannel amplifier optically connected to a digital signal acquisition system (Tucker-Davis
726 Technologies) sampling at approximately 3051.78 Hz. Line noise was removed via notch
727 filtering (60 Hz and harmonics at 120 and 180 Hz) and the resulting time series for each session
728 was visually inspected to exclude channels with excessive noise. Additionally, time segments
729 with epileptiform activity were excluded. The time series data were then common-average
730 referenced (CAR) to included electrodes either across an electrode's row in a 16x16 channel grid
731 or across the entire grid depending on the technical specifications of the amplifier used for a
732 given participant.
733
734      **High-gamma extraction.** The analytic amplitude of the high-gamma (HG; 70-150Hz)
735 frequency band was extracted by averaging across eight logarithmically-spaced bands with the
736 Hilbert transform as described elsewhere(8, 112). The HG signal was down-sampled to 400 Hz,
737 providing temporal resolution to observe latency effects on the order of <10ms (the spacing of
738 the VOTs of among the six experimental stimuli).
739
740      **Trial alignment and extraction.** Trial epochs were defined as 500ms before to 1000ms
741 after each stimulus onset. Trials were excluded for all channels if the epoch window contained
742 any time segments that had been marked for exclusion during artifact rejection. The HG signal
743 for each trial was z-scored based on the mean and standard deviation of a baseline window from
744 500ms to 200ms before stimulus onset. A 50ms moving average boxcar filter was applied to the
745 HG time series for each trial.
746
747      **Local field potential extraction.** Data for analyses of auditory evoked local field
748 potentials consisted of the same raw voltage fluctuations (local field potential), preprocessed
749 with identical notch filtering, CAR, artifact/channel rejection, and down-sampling (to 400 Hz).
750 Trial epochs (500ms before to 1000ms after each stimulus onset) were not z-scored.
751
752 **Electrode selection.**

753 **Speech-responsive electrodes.** An electrode was included in our analyses if (1) it was
754 anatomically located on the lateral temporal lobe (either superior or middle temporal gyrus), and
755 (2) the electrode's grand mean HG (across all trials and timepoints during a window 100-300ms
756 after stimulus onset) exceeded one standard deviation of the baseline window's HG activity.
757 Across all seven participants, 346 electrodes met these criteria (*speech-responsive electrodes*;
758 **Supplementary File 1; Figure 1-figure supplement 2**).
759
760 **Peak neural response.** The timepoint at which each speech-responsive electrode reached
761 its maximum HG amplitude (averaged across all trials, irrespective of condition) was identified
762 as that electrode's peak, which was used in the subsequent peak encoding analyses. Because we
763 were focused on auditory-evoked activity in the temporal lobe, the search for an electrode's peak
764 was constrained between 0 and 500ms after stimulus onset. Electrode size in **Figure 1D** and
765 **Figure 1-figure supplement 2** corresponds to this peak HG amplitude for each speech-
766 responsive electrode.
767
768 **VOT-sensitive electrodes.** To identify electrodes where the peak response depended on
769 stimulus VOT (*VOT-sensitive electrodes*), we computed the nonparametric correlation
770 coefficient (Spearman's $\rho$) across trials between VOT and peak HG amplitude. Because
771 nonparametric (rank-based) correlation analysis measures the monotonicity of the relationship
772 between two variables, it represents an unbiased ("model-free") indicator of amplitude-based
773 VOT encoding, whether the underlying monotonic relationship is categorical, linear, or follows
774 some other monotonic function (Bishara & Hittner, 2012). This procedure identified 49 VOT-
775 sensitive electrodes across all seven participants ($p < 0.05$; **Figure 2A** and **Figure 1-figure**
776 **supplement 2**; **Supplementary File 1**). Electrode color in **Figure 1D** and **Figure 1-figure**
777 **supplement 2** corresponds to the correlation coefficient at each electrode's peak (min/max $\rho$ =
778 $\pm 0.35$), thresholded such that all speech-responsive electrodes with non-significant ($p > 0.05$)
779 correlation coefficients appear as white.
780
781 This set of VOT-sensitive sites was then divided into two sub-populations based on the sign of
782 each electrode's correlation coefficient ($\rho$): voiced-selective (V+) electrodes (n = 33) had
783 significant $\rho < 0$, indicating that shorter (more /*b*/-like; voiced) VOTs elicited stronger peak HG
784 responses; voiceless-selective (V-) electrodes (n = 16) had significant $\rho > 0$, indicating that
785 longer (more /*p*/-like; voiceless) VOTs elicited stronger peak HG responses.
786
787 Across VOT-sensitive electrodes, the mean peak occurred 198.8ms after stimulus onset (SD =
788 42.3ms). The semi-transparent grey boxes in **Figures 1E** and **2B/C** illustrate this peak window
789 (mean peak $\pm$ 1 SD).
790
791 **Analysis of VOT-sensitive electrodes.**
792 **Encoding of voicing category.** Electrodes that exhibit a monotonic relationship between
793 VOT and peak HG amplitude should also be likely to exhibit a categorial distinction between
794 shorter (voiced) and longer (voiceless) VOTs. We conducted two analyses that confirmed this
795 expectation. In each analysis, we computed a nonparametric test statistic describing the
796 discriminability of responses to voiced vs. voiceless stimuli at each electrode's peak (*z*-statistic
797 of Mann-Whitney rank-sum test) and then tested whether the population of test statistics for each
798 group of electrodes (V- and V+) differed reliably from zero (Wilcoxon signed-rank tests). In the

799    first analysis, voicing category was defined based on the psychophysically determined category
800    boundary (voiced: 0-20ms VOTs; voiceless: 30-50ms VOTs), which allowed us to include all
801    VOT-sensitive electrodes (n = 49) in the analysis, including electrodes from participants who did
802    not complete the behavioral task (3/7 participants).

803

804    In the second analysis, a trial's voicing category was determined based on the actual behavioral
805    response recorded for each trial (irrespective of VOT), so this analysis was not dependent on the
806    assumption that the VOT continuum can be divided into two categories based on the average
807    boundary calculated across participants. This analysis examined the subset of trials with
808    behavioral responses and the subset of VOT-sensitive electrodes found in the four participants
809    with behavioral data (n = 27; 12 V- electrodes, 15 V+ electrodes) (**Supplementary File 1**).

810

811    Given the strong correspondence between the categorically defined VOT stimulus ranges (0-
812    20ms vs. 30-50ms VOTs) and identification behavior (e.g., **Figure 1C**), the agreement between
813    these results was expected.

814

815    Significance bars for the two example STG electrodes in one participant (e1 and e2; **Figure 1E**)
816    we computed to illustrate the temporal dynamics of category selectivity. In these electrodes, we
817    conducted the test of between-category encoding (Mann-Whitney rank-sum test; first analysis) at
818    every timepoint during the trial epoch (in addition to the electrodes' peaks). Bars plotted for each
819    electrode in **Figure 1E** begin at the first timepoint after stimulus onset where the significance
820    level reached $p < 0.005$ and ends at the first point thereafter where significance fails to reach that
821    threshold (e1: 140 to 685ms post onset; e2: 65 to 660ms post onset).

822

823    **Encoding of VOT within voicing categories.** Because VOT-sensitive electrodes were
824    identified via nonparametric correlation analysis (Spearman's $\rho$) across all VOTs, the monotonic
825    relationship between VOT and peak HG amplitude at these sites could be driven by the observed
826    phonetic (between-category) encoding of voicing without any robust sub-phonetic (within-
827    category) encoding of VOT. To assess sub-phonetic encoding of VOT in the peak response
828    amplitude of VOT-sensitive electrodes, we computed the rank-based correlation (Spearman's $\rho$)
829    between VOT and HG amplitude at each electrode's peak separately for trials in each voicing
830    category (0-20ms vs. 30-50ms VOTs). The statistical reliability of within-category encoding was
831    summarized by computing a test-statistic ($t$) for every correlation coefficient ($\rho_{0\text{-}20}$ and $\rho_{30\text{-}50}$ for
832    each VOT-sensitive electrode) as follows:

$$t = \frac{\rho\sqrt{n-2}}{\sqrt{1-\rho^2}}$$

833

834    where $n$ is the number of trials with VOTs in a given voicing category. The resulting set
835    of test statistics (one per voicing category per VOT-sensitive electrode) served as the basis for
836    the following analyses of peak within-category encoding.

837

838    For each group of electrodes (V- and V+), we tested whether the encoding of VOT within each
839    voicing category differed reliably from 0 (Wilcoxon signed-rank tests). We also conducted a
840    Wilcoxon signed-rank test for each electrode group that compared the within-category
841    correlation $t$-statistics for voiceless and voiced categories.

842

843    The above tests addressed the encoding properties of one electrode group at a time (either V- or
844    V+ electrodes). Finally, a pair of Wilcoxon signed-rank tests combined across the full set of
845    VOT-sensitive electrodes (n = 49) to summarize the within-category VOT encoding results
846    within electrodes' (1) preferred and (2) non-preferred categories. In order to conduct this
847    "omnibus" test, we multiplied the correlation *t*-statistics for all V+ electrodes (for tests within
848    each category) by -1. This simple transformation had the consequence of ensuring that positive
849    correlation statistics always indicate stronger peak HG responses to VOTs that were closer to the
850    endpoint of an electrode's preferred category.
851
852            **Visualizations of within-category VOT encoding.** To visualize the pattern of within-
853    category encoding of VOT in the peak HG amplitude of V- and V+ electrodes, we computed a
854    normalized measure of the peak response amplitude to each VOT stimulus for each VOT-
855    sensitive electrode. **Figures 2B** and **2C** show the full time series of the average (± SE) evoked
856    responses of V- and V+ electrodes to all six VOT stimuli. To show encoding patterns across
857    electrodes with different peak amplitudes, each electrode's activity was normalized by its peak
858    HG (grand mean across all VOTs). **Figure 2D** shows the amplitude of the average response
859    evoked by a given VOT at a given electrode's peak relative to the average response evoked by
860    the other VOT stimuli, or *peak HG (% of max),* averaged across electrodes in each group (V-,
861    **left**; V+, **right**) and participants (± SE). For each electrode, the mean HG amplitude evoked by
862    each VOT at the peak was scaled and normalized by subtracting the minimum across all VOTs
863    and dividing by the maximum across all VOTs after scaling.
864
865            **Neural response latency.** The normalized HG responses used for **Figures 2B/C** were
866    also used for the analysis of onset latency effects (**Figure 3**): *HG (normalized)* (**Figures 2B/C**)
867    and *HG (% of peak)* (**Figure 3A**) are computationally equivalent. Neural response onset latency
868    for an electrode was defined as the first timepoint at which its average response to a given VOT
869    stimulus exceeded 50% of its peak HG (based on the peak of the grand average response across
870    all VOTs). A bootstrapping with resampling procedure was employed to estimate the onset
871    latencies of responses to different VOTs at each electrode and to assess any possible relationship
872    between onset latency and VOT. During each sampling step in this procedure (1000 bootstrap
873    samples), we computed the average time series of the normalized HG response to each VOT, the
874    onset latency for the response to each VOT, and the nonparametric correlation (Spearman's $\rho$)
875    between onset latency and VOT. Wilcoxon signed-rank tests asked whether the population of
876    bootstrapped correlation coefficient estimates for each electrode group reliably differed from
877    zero. A Mann-Whitney rank-sum test compared the VOT-dependency of response onset latency
878    between electrode groups. Color-coded horizontal bars below the neural data in **Figure 3A** show
879    onset latency estimates (mean ± bootstrap standard error) for responses to each VOT at two
880    example electrodes. All electrodes were included in the analyses, but the bootstrapped
881    correlation coefficient estimates for two V+ electrodes that were outliers (>3 SDs from median)
882    were excluded from the visualized range of the box-plot's whiskers in **Figure 3B**.
883
884    **Population-based neural classification.** For each participant, we trained a set of multivariate
885    pattern classifiers (linear discriminant analysis with leave-one-out cross validation) to predict
886    trial-by-trial voicing category (/*b*/: 0-20ms VOTs vs. /*p*/: 30-50ms VOTs) using HG activity
887    across all speech-responsive electrodes on the temporal lobe during a time window around the
888    peak neural response. The peak window was defined as beginning 150ms and ending 250ms

889   after stimulus onset, selected based on the average and standard deviation of the peaks across all
890   VOT-sensitive electrodes. We created four separate classifiers for each participant that allowed
891   us to evaluate the contribution of amplitude and temporal structure to voicing category encoding
892   (**Figure 1F**).
893
894   To corrupt the reliability of any spatially-localized amplitude information about whether the
895   VOT stimulus presented to a participant on a given trial was a /*b*/ or a /*p*/, the neural responses at
896   every electrode on every trial were normalized so that the average response to a /*b*/ and the
897   average response to a /*p*/ reached the same amplitude at each electrode's peak. Specifically, for
898   each electrode, we found its peak (timepoint where the grand average HG time series across all
899   trials reached its maximum), calculated the mean HG amplitude across all trials for VOTs within
900   each category at that peak, and divided the HG values for every timepoint in a trial's time series
901   by the peak HG amplitude for that trial's category. This amplitude normalization procedure
902   forces the average amplitude of the neural response across all trials of /*b*/ and of /*p*/ to be equal at
903   each electrode's peak, while still allowing for variation in the amplitude of any individual trial at
904   the peak.
905
906   To corrupt the reliability of any timing information during the peak response window about
907   whether the VOT stimulus presented to a participant on a given trial was a /*b*/ or a /*p*/, the timing
908   of the neural response on every trial (across all electrodes) was randomly shifted in time so that
909   the trial could begin up to 50ms before or after the true start of the trial. Specifically, for each
910   trial, a jitter value was drawn from a discrete (integer) uniform random distribution ranging
911   between -20 to 20 (inclusive range) ECoG time samples (at 400 Hz, this corresponds to ±50ms,
912   with a mean jitter of 0ms), and the HG time series for all electrodes on that trial was moved
913   backward or forward in time by the number of samples dictated by the trial's jitter value. This
914   temporal jittering procedure has the effect of changing whether the peak response window for a
915   given trial is actually drawn from 100-200ms after stimulus onset, 200-300ms after stimulus
916   onset, or some other window in between.
917
918   Crucially, this procedure will misalign any reliable, category-dependent differences in peak
919   timing or temporal dynamics within individual electrodes or temporal patterns or relationships
920   that exist across distributed electrodes. For instance, the peak window overlaps with a window
921   during which past work examining intracranial auditory evoked local field potentials found
922   evidence of waveform shape differences between responses of single electrodes to voiced and
923   voiceless stimuli (single- vs. double-peaked responses (see, e.g., Fig. 10 of 43). If similar
924   temporal differences in waveform shape existed in the present high-gamma data, the temporal
925   jittering procedure would detect a contribution of temporal information to decoding. Moreover,
926   to the extent that the peak of a trial's evoked high-gamma response occurs during or close to the
927   peak window (either within one electrode ["local" temporal code] or across multiple electrodes
928   in the same participant ["ensemble" temporal code]), the temporal jittering procedure would
929   disrupt the reliability of this information to reveal the contribution of peak latency information to
930   decoding accuracy. On the other hand, if the peak responses to stimuli from distinct voicing
931   categories differ in the amplitude of the HG response at VOT-sensitive cortical sites, and if these
932   differences persist throughout much of the peak window, then this temporal jittering procedure is
933   unlikely to prevent the classifier from learning such differences.
934

935 For each participant, we trained one classifier where neither amplitude nor timing information
936 were corrupted (+Amplitude/+Timing), one where only timing information was corrupted
937 (+Amplitude/-Timing), one where only amplitude information was corrupted (-
938 Amplitude/+Timing), and one where both were corrupted (-Amplitude/-Timing; here, amplitude
939 normalization preceded temporal jittering). With each of these datasets, we then performed
940 dimensionality reduction to minimize overfitting using spatiotemporal principal component
941 analysis on the ECoG data for every electrode and all timepoints within the peak window
942 (retaining PCs accounting for 90% of the variance across trials of all VOTs). Finally, training
943 and testing of the linear discriminant analysis classifiers were conducted iteratively, holding out
944 a single trial, training a classifier to predict voicing category using all other trials, and then
945 predicting the voicing category of the held-out trial. For each participant and for each classifier,
946 accuracy was the proportion of held-out trials that were correctly labeled. Wilcoxon signed-rank
947 tests assessed and compared accuracy levels (across participants) achieved by the different
948 models.
949
950 **Computational neural network model.**
951        **Overview of architecture and dynamics.** A simple five-node, localist neural network
952 (**Figure 2E**) was hand-connected to illustrate how time-dependent properties of neuronal units
953 and their interactions can transform a temporal cue into a spatial code (responses of different
954 amplitudes to different VOTs at distinct model nodes). A gap detector received excitatory input
955 from both a burst detector and voicing detector, as well as input from an inhibitory node that
956 only received excitatory input from the burst detector. This represented an implementation of a
957 slow inhibitory postsynaptic potential (slow IPSP) circuit(51, 52, 60, 61). A coincidence detector
958 received excitatory input from the burst and voicing detectors.
959
960        **Network Connectivity.** Weights between units in this sparsely connected, feedforward
961 network were set according to a minimalist approach. All excitatory connections from the burst
962 detector (to the inhibitory node, the gap detector, and the coincidence detector) had identical
963 weights. All excitatory connections from the voicing detector (to the gap detector and the
964 coincidence detector) had identical weights (stronger than from burst detector). **Figure 2-figure**
965 **supplement 1** indicates all nonzero connection weights between the network's nodes, as
966 illustrated in **Figure 2E**.
967
968        **Leaky-integrator dynamics.** At the start of the model simulations, prior to the onset of
969 any stimulus ($t = 1$), the activation level $a_i(t)$ of each node $i$ was set to its resting level ($\rho_i$).
970 Simulations ran for 100 cycles, with 1 cycle corresponding to 10ms. On each subsequent cycle
971 ($t \in [2,100]$), activation levels of every node in the model were updated iteratively in two steps,
972 as described in the following algorithm:
973     (1) **Decay:** For every node $i$ with prior activation level $a_i(t-1)$ that differs from $\rho_i$, $a_i(t)$
974         decays towards $\rho_i$ by its decay rate ($\lambda_i$) without overshooting $\rho_i$.
975     (2) **Sum Inputs:** For every node $i$, the total excitatory and inhibitory inputs are summed.
976         This includes both model-external (clamped) inputs (i.e., from stimuli presented to the
977         model) on the current cycle $t$ and model-internal inputs from other nodes based on their
978         activation level on the prior cycle $a_j(t-1)$. Inputs from a presynaptic node $j$ can only
979         affect the postsynaptic node $i$ if its prior activation $a_j(t-1)$ exceeds the presynaptic
980         node's propagation threshold ($\theta_j$). Summation of model-internal inputs within $i$ is

981         weighted by the connection weights from the various presynaptic nodes (**Figure 2-figure**
982         **supplement 1**): $\sum_j w_{ji} * a_j(t-1)$. The new activation level $a_i(t)$ is bounded by the
983         node's minimum ($m_i$) and maximum ($M_i$) activation levels, irrespective of the magnitude
984         of the net effect of the inputs to a node.

985

986    All activation parameters for all nodes are listed in **Supplementary File 2**. Minimum, maximum,
987    and resting activation levels were identical across all units. Decay rates and propagation
988    thresholds were identical across the burst and voicing detectors and the inhibitory node. The
989    integrator units (gap and coincidence detectors) decayed more slowly than the other units, which
990    could only affect other model nodes during one cycle. Activation levels in the coincidence
991    detector had to reach a higher level (propagation threshold) to produce model outputs than in the
992    gap detector, a difference which allowed the gap detector to register the fast suprathreshold
993    response characteristic of slow IPSP circuits and allowed the coincidence detector to register a
994    coincidence only when both burst and voicing were detected simultaneously or at a short lag.

995

996         **Model inputs.** Two inputs were clamped onto the model in each simulation, representing
997    the onset of the burst and of voicing (**Figure 1A**). The voicing input was only clamped onto the
998    voicing detector at the onset of voicing. **Supplementary File 3** illustrates vectors describing
999    each of the simulated VOT inputs.

1000

1001         **Sensitivity of model dynamics to variations in hand-tuned model parameters.**
1002    Although most of the parameters of the model are theoretically uninteresting and were set to
1003    default levels (see **Supplementary File 2**), analysis of parameter robustness for the model
1004    revealed four primary sensitivities based on the relative values set for certain specific parameters.
1005    (1) and (2) below involve the propagation thresholds [$\theta$] of the temporal integrator units (***GAP***,
1006    ***COINC.***), which allow the model to achieve gap and coincidence detection. (3) and (4) below
1007    involve the rate of decay of activation [$\lambda$] of the temporal integrator units, which dictate where
1008    along the VOT continuum the boundary between voicing categories lies.
1009       (1) **Propagation threshold [$\theta$] of coincidence detector unit (*COINC.*):** In our model,
1010           coincidence detection is achieved by preventing the coincidence detector (***COINC.***) from
1011           propagating an output in response to the burst until the voicing has arrived (hence
1012           responding with a higher-than-minimum peak amplitude only when the voicing is
1013           coincident with or arrives shortly after the burst). Thus, the propagation threshold for
1014           ***COINC.*** ($\theta_{Coinc.}$) must be <u>greater than</u> the connection weight from the burst-detector to
1015           ***COINC.*** ($w_{Burst \rightarrow Coinc.}$).
1016       (2) **Propagation threshold [$\theta$] of gap detector unit (*GAP*):** On the other hand, the
1017           propagation threshold for the gap detector [***GAP***] ($\theta_{Gap}$) must be <u>less than</u> the connection
1018           weight from the burst-detector to ***GAP*** ($w_{Burst \rightarrow Gap}$) to register the fast suprathreshold
1019           response characteristic of slow IPSP circuits.

1020

1021    The primary factor affecting the location of the boundary between voiced (short VOTs) and
1022    voiceless (long VOTs) categories is the time-dependent rate of decay of postsynaptic potentials
1023    in ***GAP*** and ***COINC.*** towards the unit's resting activation level.
1024       (3) **Rate of decay of activation [$\lambda$] in *COINC.* in comparison to connection weights from**
1025           **inputs to *COINC.*:** For ***COINC.***, the boundary is the VOT value after which there is no
1026           longer any additional boost to its peak amplitude from the initial burst, and this requires

1027        the decay rate of ***COINC.*** ($\lambda_{Coinc.}$) and the connection weight from the burst-detector to
1028        ***COINC.*** ($w_{Burst \rightarrow Coinc.}$) to be in balance. Increasing $\lambda_{Coinc.}$ or decreasing $w_{Burst \rightarrow Coinc.}$
1029        (independently) will move the boundary earlier in time.
1030   (4) **Rate of decay of activation [$\lambda$] in *GAP* in comparison to connection weights from**
1031        **inputs to *GAP*:** Similarly, for ***GAP***, the category boundary is the VOT value before which
1032        the remaining influence of the initial inhibition is still so strong that the arrival of voicing
1033        input cannot exceed $\theta_{Gap}$. Increasing $\lambda_{Gap}$, decreasing $w_{Inhib. \rightarrow Gap}$, or increasing
1034        $w_{Voicing \rightarrow Gap}$ (independently) would each move the boundary earlier in time. All three of
1035        these parameters are in balance in these hand-tuned parameter settings.
1036
1037 It is critical to note that, for all of these cases where the hand-tuned parameter settings are in
1038 balance, the balance is required for the model to achieve gap and coincidence detection and/or to
1039 determine the position of the VOT boundary between categories. This was all the model was
1040 designed to do. No parameters were hand-tuned to achieve the other response properties (e.g.,
1041 asymmetric within-category encoding, onset latency dynamics).
1042
1043 **Analysis of auditory evoked local field potentials.**
1044      **Identification of key LFP peaks.** We identified 3 peaks of the grand mean auditory
1045 evoked local field potential (AEP), which were consistent with AEP peaks previously described
1046 in the literature(41, 42): $P_\alpha$ (positive deflection approximately 75-100 ms after stimulus onset),
1047 $N_\alpha$ (negative deflection approximately 100-150 ms after stimulus onset), and $P_\beta$ (positive
1048 deflection approximately 150-250 ms after stimulus onset) (see **Figure 1-figure supplements 3**
1049 **and 4**).
1050
1051      **Bootstrapping approach.** For each VOT-sensitive electrode (speech-responsive
1052 electrodes whose peak high-gamma amplitude was correlated with VOT), a bootstrapping with
1053 resampling procedure was used to estimate the latencies and amplitudes of each peak of the AEP
1054 elicited by trials from each VOT condition. During each sampling step in this procedure (1000
1055 bootstrap samples), we computed the average time series of the AEP for each VOT (**Figure 1-**
1056 **figure supplement 4**, **panels I-L**), the ECoG samples of the time series during each of three
1057 time-ranges with the maximum (for positive peaks) or minimum (for the negative peak) mean
1058 voltage values for each VOT, and six correlation coefficients (Pearson's *r* between VOT and
1059 amplitude/latency for each peak; see **Figure 1-figure supplement 4**, **panels M-T**).
1060
1061      **Details of peak-finding.** $P_\alpha$ was defined as the maximum mean voltage from 0-150 ms
1062 after stimulus onset, $N_\alpha$ was defined as the minimum mean voltage from 75-200 ms after
1063 stimulus onset, and $P_\beta$ was defined as the maximum mean voltage from 150-250 ms after
1064 stimulus onset. To aid peak detection and enforce sequential ordering of the peaks, time ranges
1065 for the latter two peaks ($N_\alpha$, $P_\beta$) were further constrained on a per-sample basis by setting the
1066 minimum bound of the search time range to be the time of the previous peak (i.e., the earliest
1067 possible times for $N_\alpha$ and $P_\beta$ were $P_\alpha$ and $N_\alpha$, respectively). For a given sample, if a peak
1068 occurred at either the earliest possible or latest possible time, it was assumed that the peak was
1069 either not prominent or did not occur during the defined time range for this electrode/VOT, so
1070 that sample was ignored in the analysis for that peak and any subsequent peaks. Because
1071 correlation coefficients for each peak were computed over just 6 VOTs in each sample, exclusion
1072 of a peak latency/amplitude value for one VOT condition resulted in exclusion of the all

conditions for that peak for that sample. Finally, if more than 50% of the bootstrap samples were excluded for a given peak in a given electrode, no samples for that electrode/peak pair were not included in the analysis (see, e.g., $P_\beta$ for e4 in **Figure 1-figure supplement 4**, **panels H/P/T**).

**Analysis of bootstrapped correlation estimates.** For each remaining VOT-sensitive electrode/peak pair, we determined whether or not the latency and/or amplitude of the peak was significantly associated with VOT by evaluating whether the 95% confidence interval (95% CI) across all included bootstrapped estimates of the correlation coefficient excluded 0 (taking the highest density interval of the bootstrapped statistics) (**Figure 1-figure supplement 3**, **panel B**). These exploratory analyses did not undergo multiple comparison correction.

**Detailed results of analysis of AEPs.** The exploratory analyses of correlations between VOT and the latency and/or amplitude of three peaks of the AEP in all VOT-sensitive electrodes revealed four overall conclusions:

1. Comparison of the AEPs evoked by different VOTs shows that there exist associations between stimulus VOT and the amplitude/temporal information in local field potential (LFP). Among electrodes that robustly encode voicing in their peak high-gamma amplitude (i.e., VOT-sensitive electrodes), these associations between VOT and LFP features are complex and highly variable (**Figure 1-figure supplements 3** and **4**).

2. Replicating prior results regarding VOT encoding by AEPs (e.g., 43), we find that some electrodes (e.g., e1 in **Figure 1-figure supplement 4**, **panels E/I**) exhibit temporal encoding of VOT in the latency of various peaks of the AEP. In some electrodes, the nature of this temporal code is straightforward (e.g., in e1, the latency of $N_\alpha$ is delayed by ~10ms for every additional 10ms of VOT duration; **Figure 1-figure supplement 4**, **panel M**), but – more often – the relationship between VOT and peak latency is less direct (**Figure 1-figure supplement 4**, **panels N-P**).

3. Among electrodes that encode VOT in their peak high-gamma amplitude, there exist many more electrodes that *do not* encode VOT in these temporal features of the AEP (**Figure 1-figure supplement 3**), supporting a prominent role for the peak high-gamma amplitude in the neural representation of voicing and of VOT.

4. Besides the timing of the various AEP peaks, there also exist many electrodes that encode VOT in the amplitude of those peaks (**Figure 1-figure supplement 3**). The encoding patterns are often visually similar to the encoding patterns observed in high-gamma (i.e., graded within the electrode's preferred voicing category; see **Figure 1-figure supplement 4**, **panels Q-S**). However, there are also many electrodes that do encode VOT in their peak high-gamma amplitude but *not* in these amplitude features of the LFP (**Figure 1-figure supplement 3**, **panel B**; compare, e.g., **Figure 1-figure supplement 4**, **panels D** vs. **H**).

**Supplementary analyses of spatial patterns of VOT effects.** Of the 49 VOT-sensitive electrodes, 76% were located posterior to the lateral extent of the transverse temporal sulcus (defined as $y \geq 6$ in MNI coordinate space based on projection of the sulcus onto the lateral STG in the left hemisphere). This is the same region that is densely populated with neural populations that are tuned for other phonetic features (e.g., manner of articulation(8, 82)). Mann-Whitney rank-sum tests showed that there was no significant difference in the localization of voiceless-selective (V-) versus voiced-selective (V+) electrodes along either the anterior-posterior axis (*y*-

1119  dimension in MNI coordinate space; $U = 342$, $z = -1.23$, $p = 0.22$) or the dorsal-ventral axis ($z$-
1120  dimension in MNI coordinate space; $U = 414$, $z = 0.29$, $p = 0.77$).
1121
1122  Although no regional patterns were visually apparent, we tested for hemispheric differences in
1123  relative prevalence of VOT-sensitive sites or in voicing category selectivity. Of the seven
1124  participants (all of whom had unilateral coverage), four had right hemisphere coverage (57%),
1125  and these four patients contributed 28 of the 49 VOT-sensitive electrodes identified in this study
1126  (57%) (see **Figure 2A** and **Figure 1-figure supplement 2**; **Supplementary File 1**). Pearson's $\chi^2$
1127  tests confirmed there was no difference in the rate of VOT-sensitive sites ($\chi^2(1) = 0.15$, $p = 0.70$)
1128  or in the proportion of VOT-sensitive sites that were selective for each category ($\chi^2(1) = 1.74$, $p$
1129  $= 0.19$) as a function of hemisphere. Thus, consistent with past ECoG work examining spatial
1130  patterns of STG encoding for other phonetic features(e.g., 82), we found no evidence that the
1131  observed spatial/amplitude code reflected any topographical organization nor any lateralized
1132  asymmetries in the encoding of VOT, although data limitations prevent us from ruling out this
1133  possibility entirely.
1134

**REFERENCES**

1. Stevens KN (2002) Toward a model for lexical access based on acoustic landmarks and distinctive features. *J Acoust Soc Am* 111(4):1872–1891.
2. Liberman AM, Cooper FS, Shankweiler DP, Studdert-Kennedy M (1967) Perception of the speech code. *Psychol Rev* 74(6):431–461.
3. Shannon R V, Zeng FG, Kamath V, Wygonski J, Ekelid M (1995) Speech recognition with primarily temporal cues. *Science* 270(5234):303–4.
4. Rosen S (1992) Temporal information in speech: acoustic, auditory and linguistic aspects. *Philos Trans R Soc Lond B Biol Sci* 336(1278):367–73.
5. Klatt DH (1976) Linguistic uses of segmental duration in English: Acoustic and perceptual evidence. *J Acoust Soc Am* 59(5):1208–1221.
6. Liberman AM, Delattre PC, Cooper FS (1958) Some Cues for the Distinction Between Voiced and Voiceless Stops in Initial Position. *Lang Speech* 1(3):153–167.
7. Lisker L, Abramson AS (1964) A cross-language study of voicing in initial stops: Acoustical measurements. *Word J Int Linguist Assoc* 20(3):384–422.
8. Mesgarani N, Cheung C, Johnson K, Chang EF (2014) Phonetic Feature Encoding in Human Superior Temporal Gyrus. *Science (80- )* 343(6174):1006–1010.
9. Miller JL, Green KP, Reeves A (1986) Speaking Rate and Segments: A Look at the Relation between Speech Production and Speech Perception for the Voicing Contrast. *Phonetica* 43(1–3):106–115.
10. Kessinger RH, Blumstein SE (1997) Effects of speaking rate on voice-onset time in Thai, French, and English. *J Phon* 25(2):143–168.
11. Klatt DH (1975) Voice onset time, frication, and aspiration in word-initial consonant clusters. *J Speech Hear Res* 18:686–706.
12. Lisker L, Abramson AS (1967) Some effects of context on voice onset time in English stops. *Lang Speech* 10(1):1–28.
13. Allen JS, Miller JL, DeSteno D (2003) Individual talker differences in voice-onset-time. *J Acoust Soc Am* 113(1):544.
14. Flege JE, Eefting W (1986) Linguistic and Developmental Effects on the Production and Perception of Stop Consonants. *Phonetica* 43(4):155–171.
15. Fox NP, Reilly M, Blumstein SE (2015) Phonological neighborhood competition affects spoken word production irrespective of sentential context. *J Mem Lang* 83:97–117.
16. Miller JL, Volaitis LE (1989) Effect of speaking rate on the perceptual structure of a phonetic category. *Percept Psychophys* 46(6):505–512.
17. Clayards MA, Tanenhaus MK, Aslin RN, Jacobs RA (2008) Perception of speech reflects optimal use of probabilistic speech cues. *Cognition* 108(3):804–809.
18. Kleinschmidt DF, Jaeger TF (2015) Robust speech perception: Recognize the familiar, generalize to the similar, and adapt to the novel. *Psychol Rev* 122(2). doi:10.1037/a0038695.
19. McMurray B, Jongman A (2011) What information is necessary for speech categorization? Harnessing variability in the speech signal by integrating cues computed relative to expectations. *Psychol Rev* 118(2):219–246.
20. Toscano JC, McMurray B (2010) Cue integration with categories: Weighting acoustic cues in speech using unsupervised learning and distributional statistics. *Cogn Sci* 34(3):434–464.

1191 21. Fox NP, Blumstein SE (2016) Top-down effects of syntactic sentential context on
1192     phonetic processing. *J Exp Psychol Hum Percept Perform* 42(5):730–741.
1193 22. Kuhl PK (1991) Human adults and human infants show a "perceptual magnet effect" for
1194     the prototypes of speech categories, monkeys do not. *Percept Psychophys* 50(2):93–107.
1195 23. Carney AE, Widin GP, Viemeister NF (1977) Noncategorical perception of stop
1196     consonants differing in VOT. *J Acoust Soc Am* 62(4):961–970.
1197 24. Pisoni DB, Tash J (1974) Reaction times to comparisons within and across phonetic
1198     categories. *Percept Psychophys* 15(2):285–290.
1199 25. Massaro DW, Cohen MM (1983) Categorical or continuous speech perception: A new
1200     test. *Speech Commun* 2:15–35.
1201 26. Andruski JE, Blumstein SE, Burton MW (1994) The effect of subphonetic differences on
1202     lexical access. *Cognition* 52(3):163–187.
1203 27. McMurray B, Tanenhaus MK, Aslin RN (2002) Gradient effects of within-category
1204     phonetic variation on lexical access. *Cognition* 86(2):B33–B42.
1205 28. Schouten B, Gerrits E, van Hessen A (2003) The end of categorical perception as we
1206     know it. *Speech Commun* 41(1):71–80.
1207 29. Klatt DH (1980) Software for a cascade/parallel formant synthesizer. *J Acoust Soc Am*
1208     67(3):971–995.
1209 30. Liberman AM, Harris KS, Hoffman HS, Griffith BC (1957) The discrimination of speech
1210     sounds within and across phoneme boundaries. *J Exp Psychol* 54(5):358–368.
1211 31. Liberman AM, Harris KS, Kinney JA, Lane H (1961) The discrimination of the relative
1212     onset time of the components of certain speech and non-speech patterns. *J Exp Psychol*
1213     61:379–388.
1214 32. Kronrod Y, Coppess E, Feldman NH (2016) A unified account of categorical effects in
1215     phonetic perception. *Psychon Bull Rev* 23(6):1681–1712.
1216 33. Chang EF (2015) Towards Large-Scale, Human-Based, Mesoscopic Neurotechnologies.
1217     *Neuron* 86(1):68–78.
1218 34. Crone N, et al. (2001) Induced electrocorticographic gamma activity during auditory
1219     perception. *Clin Neurophysiol* 112:565–582.
1220 35. Steinschneider M, Fishman YI, Arezzo JC (2008) Spectrotemporal Analysis of Evoked
1221     and Induced Electroencephalographic Responses in Primary Auditory Cortex (A1) of the
1222     Awake Monkey. *Cereb Cortex* 18(3):610–625.
1223 36. Ray S, Maunsell JHR (2011) Different Origins of Gamma Rhythm and High-Gamma
1224     Activity in Macaque Visual Cortex. *PLoS Biol* 9(4):e1000610.
1225 37. Steinschneider M, Volkov IO, Noh MD, Garell PC, Howard MA (1999) Temporal
1226     encoding of the voice onset time phonetic parameter by field potentials recorded directly
1227     from human auditory cortex. *J Neurophysiol* 82(5):2346–2357.
1228 38. Steinschneider M, Nourski K V., Fishman YI (2013) Representation of speech in human
1229     auditory cortex: Is it special? *Hear Res* 305(1):57–73.
1230 39. Buzsáki G, Anastassiou CA, Koch C (2012) The origin of extracellular fields and currents
1231     — EEG, ECoG, LFP and spikes. *Nat Rev Neurosci* 13(6):407–420.
1232 40. Einevoll GT, Kayser C, Logothetis NK, Panzeri S (2013) Modelling and analysis of local
1233     field potentials for studying the function of cortical circuits. *Nat Rev Neurosci*
1234     14(11):770–785.
1235 41. Howard MA, et al. (2000) Auditory cortex on the human posterior superior temporal
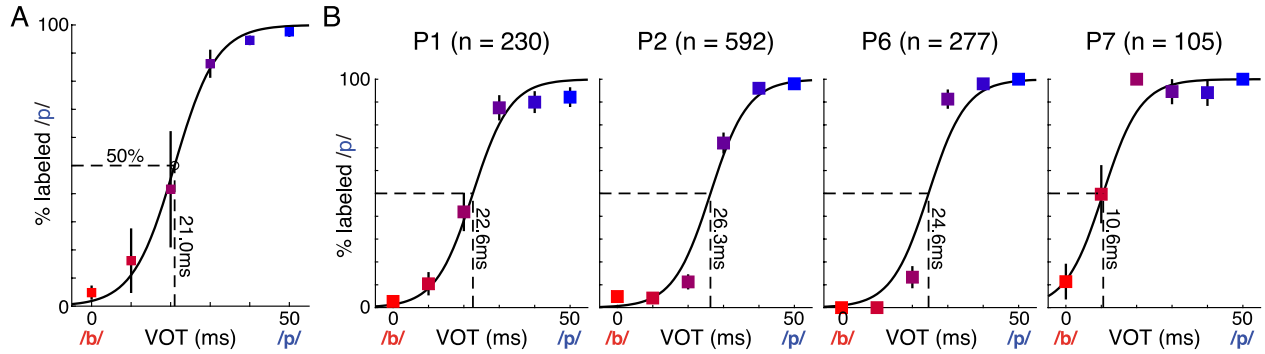1236     gyrus. *J Comp Neurol* 416(1):79–92.

1237    42.    Nourski K V, et al. (2015) Sound identification in human auditory cortex: Differential
1238           contribution of local field potentials and high gamma power as revealed by direct
1239           intracranial recordings. *Brain Lang* 148:37–50.
1240    43.    Steinschneider M, et al. (2011) Intracranial study of speech-elicited activity on the human
1241           posterolateral superior temporal gyrus. *Cereb Cortex* 21(Cv):2332–47.
1242    44.    Blumstein SE, Myers EB, Rissman J (2005) The perception of voice onset time: an fMRI
1243           investigation of phonetic category structure. *J Cogn Neurosci* 17(9):1353–66.
1244    45.    Toscano JC, Mcmurray B, Dennhardt J, Luck SJ (2010) Continuous perception and
1245           graded categorization: electrophysiological evidence for a linear relationship between the
1246           acoustic signal and perceptual encoding of speech. *Psychol Sci* 21(10):1532–1540.
1247    46.    Toscano JC, Anderson ND, Fabiani M, Gratton G, Garnsey SM (2018) The time-course of
1248           cortical responses to speech revealed by fast optical imaging. *Brain Lang* 184:32–42.
1249    47.    Frye RE, et al. (2007) Linear coding of voice onset time. *J Cogn Neurosci* 19(9):1476–
1250           1487.
1251    48.    Myers EB (2007) Dissociable effects of phonetic competition and category typicality in a
1252           phonetic categorization task: An fMRI investigation. *Neuropsychologia* 45(7):1463–1473.
1253    49.    Ferster D, Spruston N (1995) Cracking the neuronal code. *Science (80- )* 270(5237):756–
1254           757.
1255    50.    Shadlen MN, Newsome WT (1994) Noise, neural codes and cortical organization. *Curr
1256           Opin Neurobiol* 4(4):569–579.
1257    51.    Buonomano D V., Merzenich MM (1995) Temporal information transformed into a
1258           spatial code by a neural network with realistic properties. *Science (80- )*
1259           267(February):1028–1030.
1260    52.    Gao X, Wehr M (2015) A Coding Transformation for Temporally Structured Sounds
1261           within Auditory Cortical Neurons. *Neuron* 86(1):292–303.
1262    53.    Eggermont JJ (2000) Neural Responses in Primary Auditory Cortex Mimic
1263           Psychophysical, Across-Frequency-Channel, Gap-Detection Thresholds. *J Neurophysiol*
1264           84(3):1453–1463.
1265    54.    Carr CE (1993) Processing of Temporal Information in the Brain. *Annu Rev Neurosci*
1266           16(1):223–243.
1267    55.    Konishi M (2003) Coding of Auditory Space. *Annu Rev Neurosci* 26(1):31–55.
1268    56.    Rauschecker JP (2014) Is there a tape recorder in your head? How the brain stores and
1269           retrieves musical melodies. *Front Syst Neurosci* 8:149.
1270    57.    Rauschecker JP (1998) Cortical processing of complex sounds. *Curr Opin Neurobiol*
1271           8(4):516–521.
1272    58.    McClelland JL, Rumelhart DE (1981) An interactive activation model of context effects in
1273           letter perception. *Psychol Rev* 88:375–407.
1274    59.    McClelland JL, Mirman D, Bolger DJ, Khaitan P (2014) Interactive Activation and
1275           Mutual Constraint Satisfaction in Perception and Cognition. *Cogn Sci* 38(6):1139–89.
1276    60.    Douglas RJ, Martin KA (1991) A functional microcircuit for cat visual cortex. *J Physiol*
1277           440:735–69.
1278    61.    McCormick DA (1989) GABA as an inhibitory neurotransmitter in human cerebral cortex.
1279           *J Neurophysiol* 62(5):1018–27.
1280    62.    Margoliash D, Fortune ES (1992) Temporal and harmonic combination-sensitive neurons
1281           in the zebra finch's HVc. *J Neurosci* 12(11):4309–26.
1282    63.    Peña JL, Konishi M (2001) Auditory Spatial Receptive Fields Created by Multiplication.

1283    *Science (80- )* 292(5515):249–252.

64.    Peña JL, Konishi M (2002) From Postsynaptic Potentials to Spikes in the Genesis of Auditory Spatial Receptive Fields. *J Neurosci* 22(13):5652–5658.

65.    Lisker L (1986) "Voicing" in English: A catalogue of acoustic features signaling /b/ versus /p/ in trochees. *Lang Speech* 29(1):3–11.

66.    Soli SD (1983) The role of spectral cues in discrimination of voice onset time differences. *J Acoust Soc Am* 73(6):2150–2165.

67.    Stevens KN, Klatt DH (1974) Role of formant transitions in the voiced-voiceless distinction for stops. *J Acoust Soc Am* 55(3):653–659.

68.    Summerfield Q, Haggard M (1977) On the dissociation of spectral and temporal cues to the voicing distinction in initial stop consonants. *J Acoust Soc Am* 62(2):435–448.

69.    Eggermont JJ (1995) Representation of a voice onset time continuum in primary auditory cortex of the cat. *J Acoust Soc Am* 98(2):911–920.

70.    Eggermont JJ, Ponton CW (2002) The neurophysiology of auditory perception: from single units to evoked potentials. *Audiol Neurootol* 7(2):71–99.

71.    Liégeois-Chauvel C, De Graaf JB, Laguitton V, Chauvel P (1999) Specialization of left auditory cortex for speech perception in man depends on temporal coding. *Cereb Cortex* 9(5):484–496.

72.    Steinschneider M, et al. (2005) Intracortical responses in human and monkey primary auditory cortex support a temporal processing mechanism for encoding of the voice onset time phonetic. *Cereb Cortex* 15:170–186.

73.    Steinschneider M, Schroeder CE, Arezzo JC, Vaughan HG (1994) Speech-evoked activity in primary auditory cortex: effects of voice onset time. *Electroencephalogr Clin Neurophysiol* 92:30–43.

74.    Steinschneider M, Schroeder CE, Arezzo JC, Vaughan HG (1995) Physiologic correlates of the voice onset time boundary in primary auditory cortex (A1) of the awake monkey: temporal response patterns. *Brain Lang* 48(3):326–340.

75.    Steinschneider M, Fishman YI, Arezzo JC (2003) Representation of the voice onset time (VOT) speech parameter in population responses within primary auditory cortex. *J Acoust Soc Am* 114(1):307–321.

76.    Theunissen FE, Miller JP (1995) Temporal encoding in nervous systems: A rigorous definition. *J Comput Neurosci* 2(2):149–162.

77.    Engineer CT, et al. (2008) Cortical activity patterns predict speech discrimination ability. *Nat Neurosci* 11(5):603–8.

78.    Eggermont JJ (2001) Between sound and perception: reviewing the search for a neural code. *Hear Res* 157(1–2):1–42.

79.    Oxenham AJ (2018) How We Hear: The Perception and Neural Coding of Sound. *Annu Rev Psychol* 69(1):27–50.

80.    Yi HG, Leonard MK, Chang EF (2019) The Encoding of Speech Sounds in the Superior Temporal Gyrus. *Neuron* 102(6):1096–1110.

81.    Tang C, Hamilton LS, Chang EF (2017) Intonational speech prosody encoding in the human auditory cortex. *Science (80- )* 357(6353):797–801.

82.    Hamilton LS, Edwards E, Chang EF (2018) A Spatial Map of Onset and Sustained Responses to Speech in the Human Superior Temporal Gyrus. *Curr Biol* 28(12):1860-1871.e4.

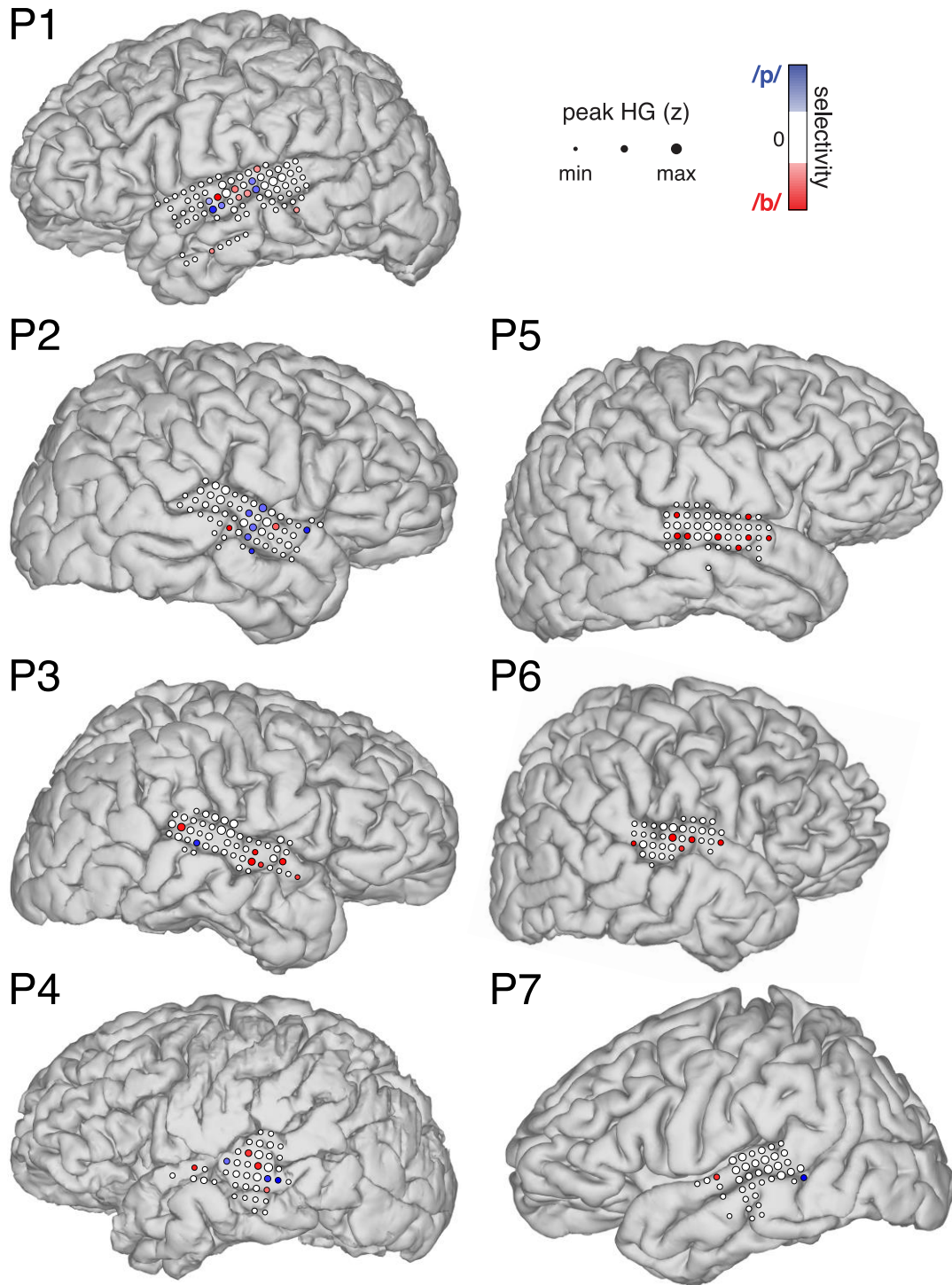83.    Oganian Y, Chang EF (2019) A speech envelope landmark for syllable encoding in human

1329    superior temporal gyrus. *Sci Adv* 5(11):eaay6279.

1330 84. McClelland JL, Elman JL (1986) The TRACE model of speech perception. *Cogn Psychol*
1331    18(1):1–86.

1332 85. Norris D, McQueen JM (2008) Shortlist B: A Bayesian model of continuous speech
1333    recognition. *Psychol Rev* 115(2):357–395.

1334 86. Norris D, McQueen JM, Cutler A (2015) Prediction, Bayesian inference and feedback in
1335    speech recognition. *Lang Cogn Neurosci*:1–15.

1336 87. Magnuson J, et al. EARSHOT: A minimal neural network model of incremental human
1337    speech recognition. doi:10.31234/OSF.IO/H7A4N.

1338 88. Damper RI (1994) Connectionist models of categorical perception of speech. *Proceedings*
1339    *of ICSIPNN 1994 International Symposium on Speech, Image Processing and Neural*
1340    *Networks* (Institute of Electrical and Electronics Engineers Inc.), pp 101–104.

1341 89. Kössl M, et al. (2014) Neural maps for target range in the auditory cortex of echolocating
1342    bats. *Curr Opin Neurobiol* 24:68–75.

1343 90. Portfors C V., Wenstrup JJ (2001) Topographical distribution of delay-tuned responses in
1344    the mustached bat inferior colliculus. *Hear Res* 151(1–2):95–105.

1345 91. Zatorre RJ, Belin P (2001) Spectral and temporal processing in human auditory cortex.
1346    *Cereb Cortex* 11(10):946–953.

1347 92. Fries P (2009) Neuronal Gamma-Band Synchronization as a Fundamental Process in
1348    Cortical Computation. *Annu Rev Neurosci* 32(1):209–224.

1349 93. Giraud A-L, Poeppel D (2012) Cortical oscillations and speech processing: emerging
1350    computational principles and operations. *Nat Neurosci* 15(4):511–517.

1351 94. Kösem A, et al. (2018) Neural Entrainment Determines the Words We Hear. *Curr Biol*
1352    28(18):2867-2875.e3.

1353 95. Peelle JE, Davis MH (2012) Neural Oscillations Carry Speech Rhythm through to
1354    Comprehension. *Front Psychol* 3:320.

1355 96. Chang EF, et al. (2010) Categorical speech representation in human superior temporal
1356    gyrus. *Nat Neurosci* 13(11):1428–32.

1357 97. Macmillan NA, Kaplan HL, Creelman CD (1977) The psychophysics of categorical
1358    perception. *Psychol Rev* 84(5):452–471.

1359 98. Lee YS, Turkeltaub P, Granger R, Raizada RDS (2012) Categorical speech processing in
1360    Broca's area: An fMRI study using multivariate pattern-based analysis. *J Neurosci*
1361    32(11):3942–3948.

1362 99. Myers EB, Blumstein SE, Walsh E, Eliassen J (2009) Inferior Frontal Regions Underlie
1363    the Perception of Phonetic Category Invariance. *Psychol Sci* 20(7):895–903.

1364 100. Evans S, Davis MH (2015) Hierarchical organization of auditory and motor
1365    representations in speech perception: evidence from searchlight similarity analysis. *Cereb*
1366    *Cortex* 25:4772–4788.

1367 101. Sohoglu E, Peelle JE, Carlyon RP, Davis MH (2012) Predictive top-down integration of
1368    prior knowledge during speech perception. *J Neurosci* 32(25):8443–53.

1369 102. Leonard MK, Baud MO, Sjerps MJ, Chang EF (2016) Perceptual restoration of masked
1370    speech in human cortex. *Nat Commun* 7:13619.

1371 103. Cope TE, et al. (2017) Evidence for causal top-down frontal contributions to predictive
1372    processes in speech perception. *Nat Commun* 8(1):2154.

1373 104. Park H, Ince RAA, Schyns PG, Thut G, Gross J (2015) Frontal Top-Down Signals
1374    Increase Coupling of Auditory Low-Frequency Oscillations to Continuous Speech in

1375        Human Listeners. *Curr Biol* 25(12):1649–1653.

1376  105.  McClelland JL, Mirman D, Holt LL (2006) Are there interactive processes in speech
1377        perception? *Trends Cogn Sci* 10(8):363–369.

1378  106.  McQueen JM, Norris D, Cutler A (2006) Are there really interactive processes in speech
1379        perception? *Trends Cogn Sci* 10(12). doi:10.1016/j.tics.2006.10.004.

1380  107.  Norris D, McQueen JM, Cutler A (2000) Merging information in speech recognition:
1381        feedback is never necessary. *Behav Brain Sci* 23(3):299–325.

1382  108.  Cho T, Ladefoged P (1999) Variation and universals in VOT: evidence from 18
1383        languages. *J Phon* 27(2):207–229.

1384  109.  DeWitt I, Rauschecker JP (2012) Phoneme and word recognition in the auditory ventral
1385        stream. *Proc Natl Acad Sci U S A* 109(8):E505-14.

1386  110.  Obleser J, Eisner F (2009) Pre-lexical abstraction of speech in the auditory cortex. *Trends*
1387        *Cogn Sci* 13(1):14–19.

1388  111.  Leonard MK, Chang EF (2014) Dynamic speech representations in the human temporal
1389        lobe. *Trends Cogn Sci* 18(9):472–479.

1390  112.  Sjerps MJ, Fox NP, Johnson K, Chang EF (2019) Speaker-normalized sound
1391        representations in the human auditory cortex. *Nat Commun* 10(1):2465.

1392  113.  Fox NP, Leonard MK, Sjerps MJ, Chang EF (2020) Transformation of a temporal speech
1393        cue to a spatial neural code in human auditory cortex. *Open Sci Framew*. Available at:
1394        https://osf.io/9y7uh/.

1395  114.  Hamilton LS, Chang DL, Lee MB, Chang EF (2017) Semi-automated Anatomical
1396        Labeling and Inter-subject Warping of High-Density Intracranial Recording Electrodes in
1397        Electrocorticography. *Front Neuroinform* 11:62.

1398  115.  Feldman NH, Griffiths TL, Morgan JL (2009) The influence of categories on perception:
1399        explaining the perceptual magnet effect as optimal statistical inference. *Psychol Rev*
1400        116(4):752–82.
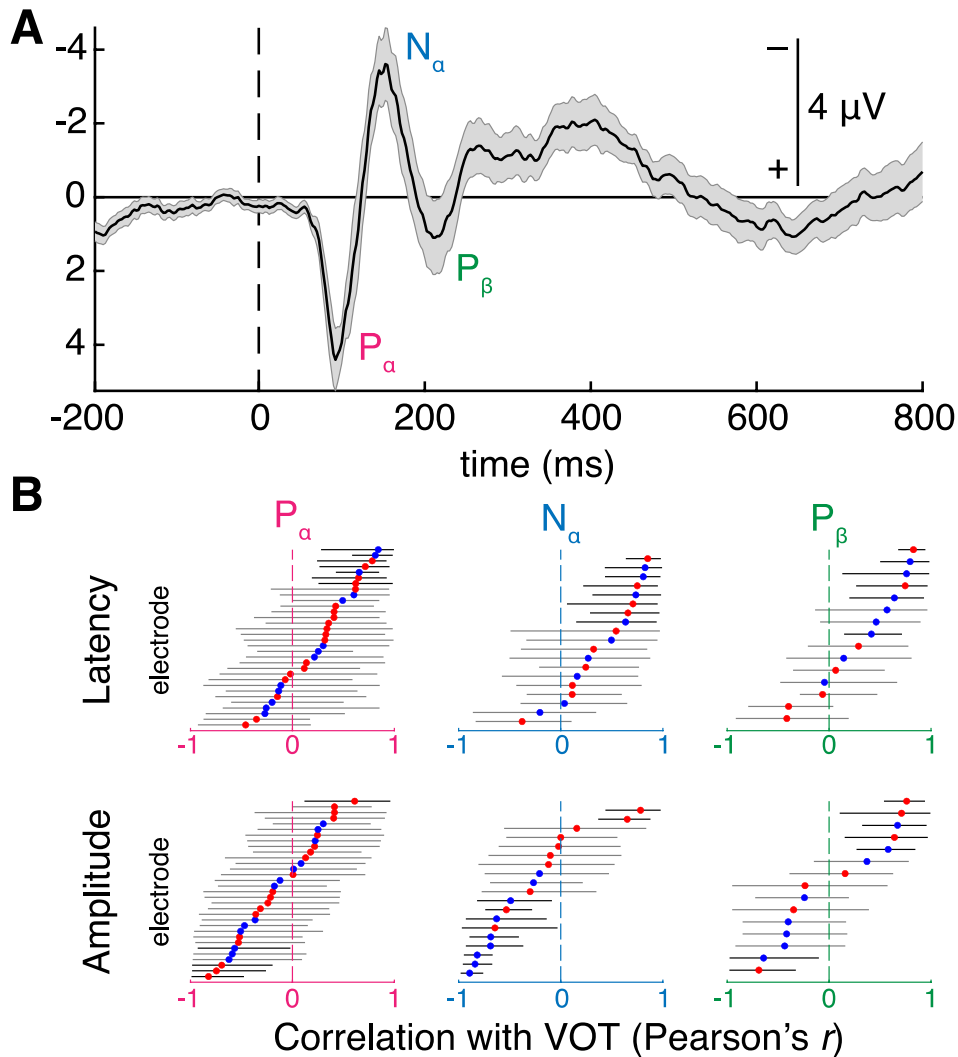
1401

1402 **FIGURE SUPPLEMENTS**

1403

1404



1405

1406 ***Figure 1-figure supplement 1. Identification behavior across all participants with behavioral data. A.***
1407 *Mean (± SE across participants; n = 4 of 7 participants) percent /pa/ responses for each voice-onset time*
1408 *(VOT) stimulus. Best-fit psychometric curve (mixed effects logistic regression) yields voicing category*
1409 *boundary at 21.0ms (50% crossover point; see **Methods** for details).* ***B.*** *Behavior (mean ± bootstrap SE)*
1410 *for each individual participant (P1, P2, P6, P7). Total trials (n) listed for each participant (see*
1411 ***Supplementary File 1***). *Best-fit psychometric curves and category boundaries were computed using the*
1412 *mixed effects logistic regression across all participants, adjusted by the random intercept fit by the model*
1413 *for each participant. Voicing category boundaries were subject-dependent, with 3 of 4 participants'*
1414 *occurring between 20-30ms. P1 is representative participant in **Figure 1C**.*
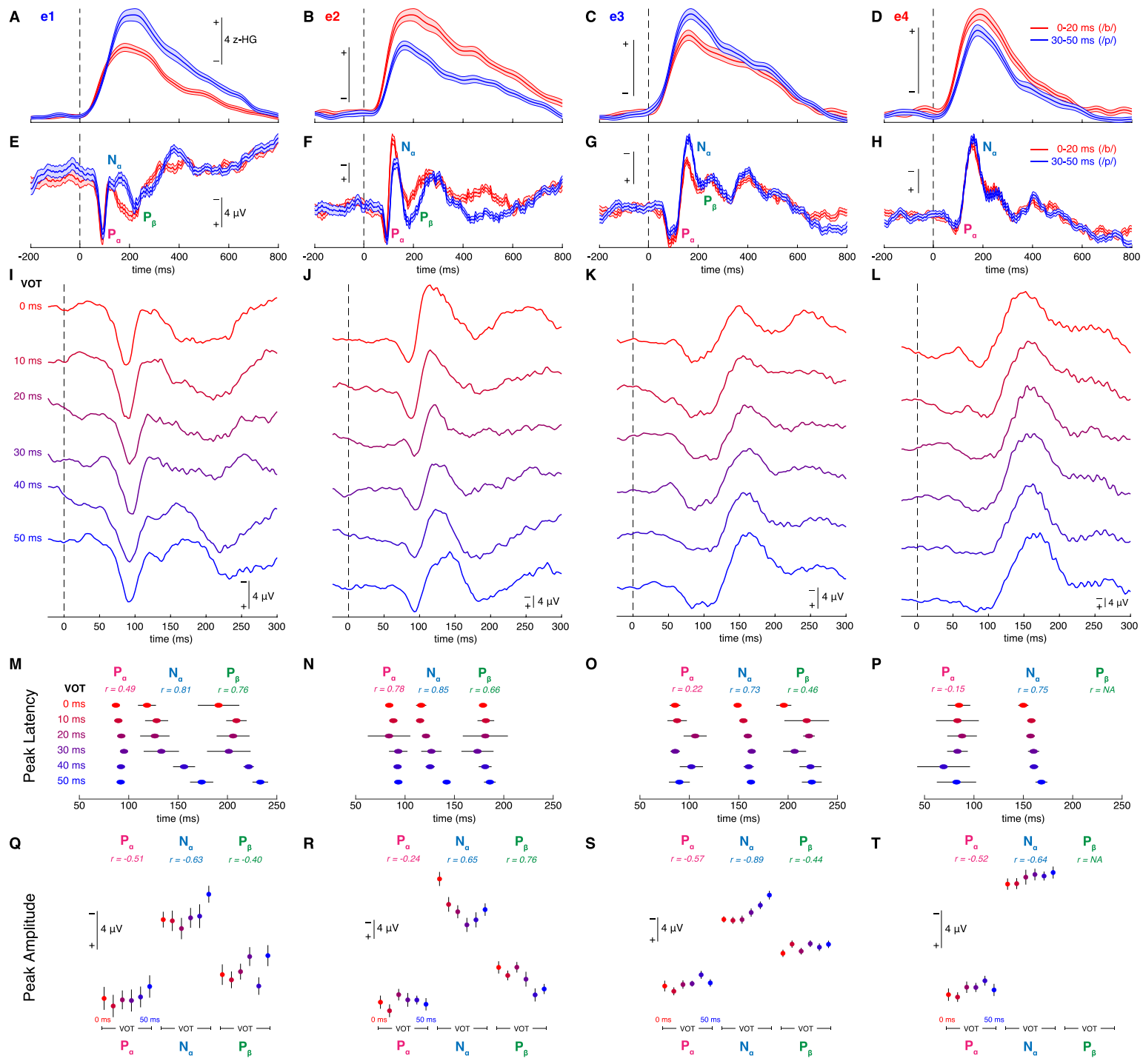
1415

1416
1417 ***Figure 1-figure supplement 2. Locations of all speech-responsive and VOT-sensitive electrodes in each***
1418 ***participant (P1-P7).*** *P1 is representative participant in **Figure 1D**. Electrode color reflects strength and*
1419 *direction of selectivity (Spearman's ρ between peak HG amplitude and VOT) at subset of VOT-sensitive*
1420 *sites (p < 0.05) for either voiceless VOTs (/p/; blue) or voiced VOTs (/b/; red). Electrode size indicates*
1421 *peak high-gamma (HG; z-scored) amplitude at all speech-responsive temporal lobe sites. Maximum and*
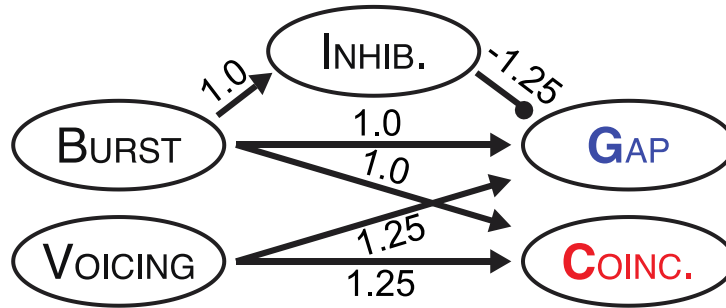1422 *minimum electrode size and selectivity was calculated per participant for visualization.*
1423

1424
1425 *Figure 1-figure supplement 3. Analysis of evoked local field potentials reveals that some electrodes that*
1426 *encode VOT in their peak high-gamma amplitude also exhibit amplitude and/or temporal response*
1427 *features that are VOT-dependent. A. Grand average auditory evoked potential (AEP) to all VOT stimuli.*
1428 *Evoked local field potentials (negative up-going) were averaged over all VOT-sensitive STG electrodes*
1429 *for one representative participant (P1) (mean ± SE, computed across electrodes). Three peaks of the AEP*
1430 *were identified for analysis: 75-100 ms (Pα), 100-150 ms (Nα), and 150-250 ms (Pβ) after stimulus onset.*
1431 ***B.*** *Correlation coefficients (Pearson's r) quantifying association between VOT and latency (**top**) or*
1432 *amplitude (**bottom**) of each peak (Pα: **left**; Nα: **middle**; Pβ: **right**) for each VOT-sensitive electrode for*
1433 *which that peak could be reliably identified (see **Figure 1-figure supplement 4** and **Methods** for details*
1434 *of this analysis). Horizontal bars represent bootstrapped estimate of correlation coefficient (mean and*
1435 *95% CI) for each electrode (blue: voiceless-selective; red: voiced-selective; electrodes sorted by mean*
1436 *correlation value). Black bars around an electrode's mean indicate that encoding of VOT by the*
1437 *designated parameter (latency or amplitude of a given peak) was significant (95% CI excluded r = 0;*
1438 *grey bars: not significant). Later peaks were reliably identified for fewer electrodes (Pα: n = 32 of 49*
1439 *electrodes; Nα: n = 19; Pβ: n = 15).*
1440

1442 *Figure 1-figure supplement 4. Complex and variable associations between VOT and*
1443 *amplitude/temporal features of auditory evoked local field potentials (AEPs) exist in responses of*
1444 *electrodes that robustly encode voicing in their peak high-gamma amplitude. A to D. Average high-*
1445 *gamma responses (± SE) to voiced (0-20ms VOTs; red) and voiceless (30-50ms VOTs; blue) stimuli in*
1446 *four representative VOT-sensitive STG electrodes, including two voiceless-selective (A: e1, C: e3) and*
1447 *two voiced-selective (B: e2, D: e4) electrodes, aligned to stimulus onset. Vertical bars indicate relative*
1448 *scaling of high-gamma (z-scored) in each panel. The two leftmost electrodes (e1, e2) correspond to e1*
1449 *and e2 in main text (e.g., Figure 1E). E to H. Average local field potentials (± SE) evoked by*
1450 *voiced/voiceless stimuli in the same four electrodes, aligned to stimulus onset. Vertical bars (negative-*
1451 *upgoing) indicate relative scaling of voltage in each panel. The three peaks of the AEP that were*
1452 *identified for analysis are labeled for each electrode ($P_\alpha$, $N_\alpha$, $P_\beta$; see Figure 1-figure supplement 3). For*
1453 *a given electrode, peaks were omitted from this analysis if they could not be reliably identified across*
1454 *bootstrapped samples of trials from all six VOT conditions (e.g., $P_\beta$ for e4). See Methods for details. I to*
1455 *L. Average local field potentials evoked by each VOT stimulus (line color) in the same four electrodes,*
1456 *aligned to stimulus onset. M to P. Mean latency (± bootstrap SE) of each AEP peak for each VOT*
1457 *stimulus for the same four electrodes. Mean bootstrapped correlation (Pearson's r) between VOT and*
1458 *peak latency shown for each peak/electrode. Q to T. Mean amplitude (± bootstrap SE) of each AEP peak*
1459 *for each VOT stimulus for the same four electrodes. Mean bootstrapped correlation (Pearson's r)*
1460 *between VOT and peak amplitude shown for each peak/electrode. Note that negative correlations are*
1461 *visually represented as rising from left to right. Correlation coefficients comprised the source data for*
1462 *summary representations in Figure 1-figure supplement 3.*
1463

1464
1465     *Figure 2-figure supplement 1. Connection weights between model nodes.*
1466

**LEGENDS FOR SUPPLEMENTARY FILES**

1468

| Participant | Hem | # trials (ECoG) | # trials (behavior) | # elecs (SR) | # elecs (VOT) | # elecs (V- / V+) |
|---|---|---|---|---|---|---|
| P1 | LH | 234 | 230 | 78 | 12 | 5 / 7 |
| P2 | RH | 625 | 592 | 56 | 8 | 6 / 2 |
| P3 | RH | 339 | 0 | 50 | 7 | 1 / 6 |
| P4 | LH | 333 | 0 | 40 | 7 | 3 / 4 |
| P5 | RH | 119 | 0 | 47 | 8 | 0 / 8 |
| P6 | RH | 305 | 277 | 36 | 5 | 0 / 5 |
| P7 | LH | 110 | 105 | 39 | 2 | 1 / 1 |

1469 ***Supplementary File 1. Table of experimental summary statistics for each participant.*** *Each participant*
1470 *had ECoG grid coverage of one hemisphere (Hem), either left (LH) or right (RH). Participants completed*
1471 *as many trials as they felt comfortable with. Number of trials per participant for ECoG analyses indicate*
1472 *trials remaining after artifact rejection. Some participants chose to listen passively to some or all blocks,*
1473 *so three participants have no trials for behavioral analyses. See* ***Methods*** *for description of inclusion*
1474 *criteria for individual trials in ECoG and behavioral analyses. A subset of speech-responsive (SR)*
1475 *electrodes on the lateral surface of the temporal lobe had a peak amplitude that was sensitive to VOT,*
1476 *selectively responding to either voiceless (V-) or voiced (V+) stimuli. See* ***Methods*** *for details on*
1477 *electrode selection.*

1478

| | | activation parameter | | | | |
|---|---|---|---|---|---|---|
| | | $m$ | $M$ | $\rho$ | $\lambda$ | $\theta$ |
| model node | Burst | -10 | 10 | 0 | 1 | 0 |
| | Voicing | -10 | 10 | 0 | 1 | 0 |
| | Inhibitor | -10 | 10 | 0 | 1 | 0 |
| | Gap | -10 | 10 | 0 | 0.25 | 0.25 |
| | Coincidence | -10 | 10 | 0 | 0.25 | 1 |

1479 ***Supplementary File 2. Table of activation parameters for each model node.*** $m$ *= minimum activation*
1480 *level.* $M$ *= maximum activation level.* $\rho$ *= resting activation level.* $\lambda$ *= decay rate.* $\theta$ *= propagation*
1481 *threshold.*

1482

| VOT | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | | | BV | | | | | | | | | | | | |
| 10 | | | B | V | | | | | | | | | | | |
| 20 | | | B | | V | | | | | | | | | | |
| 30 | | | B | | | V | | | | | | | | | |
| 40 | | | B | | | | V | | | | | | | | |
| 50 | | | B | | | | | V | | | | | | | |
| | -2 | -1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | … |
| | time post onset (cycles) | | | | | | | | | | | | | | |

1483 ***Supplementary File 3. Table illustrating timing of 6 simulated model inputs.*** *The table is sparse,*
1484 *meaning that inputs to both Burst and Voicing detector units are 0 whenever a cell is blank. Inputs are*
1485 *clamped onto either Burst or Voicing detector units (always with strength = 1) for a given simulated VOT*
1486 *stimulus during the cycles that are labeled with a B or a V.*

1487