

24 To whom correspondence should be addressed:

25 Ido Davidesco, Psychology Department, New York University, New York, NY 10003

26 Email: ido.davidesco@nyu.edu

27

28 Running Head: ECoG responses to time-compressed speech

29

30 **Abstract**

31 Human listeners understand spoken language across a variety of rates, but when speech is
32 presented three times or more faster than its usual rate, it becomes unintelligible. How the brain
33 achieves such tolerance and why speech becomes unintelligible above certain rates is still
34 unclear. We addressed these questions using electrocorticography (ECoG) recordings in 7
35 epileptic patients (two female). Patients rated the intelligibility of sentences presented at the
36 original rate (100%), speeded rates (33% or 66% of the original sentence duration) and a slowed
37 rate (150%). We then examined which parameters of the neural response covary with the
38 transition from intelligible to unintelligible speech. Specifically, we asked whether neural
39 responses: 1) track the acoustic envelope of the incoming speech; 2) “scale” with speech rate, i.e.
40 whether neural responses elicited by slowed and speeded sentences can be linearly scaled to
41 match the responses to the original sentence. Behaviorally, intelligibility was at ceiling for
42 speech rates of 66% and above, but dropped significantly for the 33% rate. At the neural level,
43 Superior Temporal Gyrus regions (STG) in close proximity to A1 (‘low-level’) tracked the
44 acoustic envelope and linearly scaled with the input across all speech rates, irrespective of
45 intelligibility. In contrast, secondary auditory areas in the STG as well as the inferior frontal
46 gyrus and angular gyrus (‘high-level’) tracked the acoustic envelope and linearly scaled with

47 input only for intelligible speech. These results help reconcile seemingly contradictory previous
48 findings and provide better understanding of how information processing unfolds along the
49 cortical auditory hierarchy.

50

51 **Keywords**

52 Electrocorticography; time-compressed speech; Speech intelligibility

53

54 **New & Noteworthy**

55 The human brain can cope with large variations in speech rate. However, when speech is
56 artificially accelerated, above a certain rate it becomes incomprehensible. This study investigated
57 how the brain achieves this tolerance to speech rate, and what might constrain our understanding
58 of speeded-up speech. Whereas in low-level auditory areas, neural responses scaled with speech
59 rate irrespective of intelligibility, high-order brain regions could only track speech as long as it
60 remained comprehensible.

61

62

63 **Introduction**

64 Human listeners understand speech over a wide range of rates. Speech remains intelligible even
65 when it is artificially slowed or accelerated up to 40% of its original duration (Dupoux and Green
66 1997; Mehler et al. 1993; Pallier et al. 1998; Sebastián-Gallés et al. 2000). However, how this
67 tolerance to temporal variability is achieved at the neural level and why spoken language
68 becomes unintelligible above certain rates is currently poorly understood.

69 Nourski et al. (2009) demonstrated that high-frequency (>70 Hz) electrocorticographic (ECoG)
70 responses recorded directly from Heschl's gyrus (A1) could track the speech envelope well
71 outside of the intelligibility range. On the other hand, Ahissar et al. (2001) reported that time
72 compression of speech beyond the intelligibility limit is associated with a sharp decrease in the
73 temporal locking of auditory magnetoencephalographic (MEG) responses to the speech
74 envelope. More recently, using functional MRI, Lerner et al. (2014) measured blood-
75 oxygenation level dependent (BOLD) responses to speeded-up and slowed-down versions of a 7-
76 minute narrative (50% to 200%). They found that for both the slowed-down and speeded-up
77 rates, linearly scaled BOLD responses matched the response to the original narrative. This linear
78 scaling of the neural responses was observed across the entire processing hierarchy, including
79 early auditory regions as well as linguistic and extra-linguistic brain areas (but note that speech
80 was always kept within the intelligibility range).

81 Although the findings described above seem to be contradictory, it is possible that they reflect
82 different stages of processing along the auditory hierarchy. In a series of studies, we have
83 demonstrated a neural hierarchy of Temporal Receptive Windows (TRWs) (Hasson et al. 2008;
84 Honey et al. 2012; Lerner et al. 2011). Analogous to the notion of a spatial receptive field, TRW
85 refers to the window of time in which information is being integrated. The TRW gradually

86 increases from early sensory areas to higher-order perceptual and cognitive areas (Lerner et al.
87 2011). Therefore, we hypothesize that the short temporal integration windows of early auditory
88 areas (e.g. A1) would enable the tracking of accelerated speech even outside of the intelligibility
89 range. In contrast, in higher order areas, the integration of information may fail at high
90 compression rates.

91 In the current study, we used ECoG recordings in seven neurosurgical patients to address the
92 question of where along the cortical processing timescale hierarchy invariance to speech rate
93 emerges. Participants were presented with a list of sentences spoken at a normal rate (100%) as
94 well as slowed-down (150% duration) and speeded-up (66% and 33%) rates. Following Nourski
95 et al. (2009) and Ahissar et al. (2001), we correlated the speech envelope with the envelope of
96 the broadband (75-200Hz) neural responses at each speech rate. Based on the Lerner et al. (2014)
97 study, we also tested the extent to which linear scaling of the neural responses elicited by
98 speeded (or slowed down) sentences match the neural responses to the original speech rate.
99 Whereas neural tracking is mostly sensitive to low-level properties of the speech signal (i.e.
100 variations in amplitude across time), linear scaling can capture more high-level properties of
101 speech processing (Lerner et al., 2014). Even though we did not have access to neural data from
102 A1, we predicted that adjacent early auditory areas along the STG would exhibit speech rate
103 invariance irrespective of intelligibility level. In contrast, areas outside of early auditory cortex,
104 further along the processing hierarchy, which integrate sounds into intelligible syllables and
105 words, would scale their neural activity with speech rate only within the intelligibility range.

106

107

108

109 **Materials and Methods**

110 *Participants*

111 Seven native speakers of English (2 female; 24-56 years old) experiencing pharmacologically
112 refractory complex partial seizures were recruited via the Comprehensive Epilepsy Center of the
113 New York University School of Medicine. Their clinical and demographic information is
114 summarized in Table 1. Patients had elected to undergo intracranial monitoring for clinical
115 purposes and provided written informed consent in accordance with New York University
116 Medical Center Institutional Review Board. Electrode placement was determined based on
117 clinical criteria without reference to this study. Patients had left-hemisphere (n=3), right-
118 hemisphere (n=3) and bilateral (n=1) electrode coverage.

119

120 **Table 1:** Demographic and Recording Characteristics of Patients.

Patient ID	Gender	Age (years)	WADA ¹ language	Implanted hemisphere	# of implanted electrodes
NY393	M	44	Left	Left	120
NY339	F	24	N/A	Left	122
NY415	F	56	Left	Left	124
NY400	M	27	Left	Bilateral	124
NY442	M	29	N/A	Right	64
NY394	M	27	Left	Right	124
NY451	M	25	N/A	Right	204

121 ¹ Also known as intracarotid sodium amobarbital procedure; used to map language localization.

122

123 *Stimuli*

124 A set of 33 spoken sentences with duration ranging between 3 and 3.3 seconds were selected
125 from the Harvard sentences corpus (IEEE 1969). All sentences were recorded by a male speaker.

126 Stimuli covered four speech rates: uncompressed (100%) speech, slowed (150%) and speeded
127 speech (33% and 66% of the duration of the corresponding uncompressed signal; See Fig. 1A).

128 Unfortunately, we could not include more intermediate speech rates because of the limited
129 testing time available with each patient. The original rate and 33% conditions were represented
130 by 33 sentences and the 66% and 150% conditions – by 25 sentences that were randomly
131 selected from the set of 33. Sentences were presented consecutively, in pseudorandom order,
132 until each sentence had been presented twice.

133 To control for sentence duration, we generated concatenated (C) sentences which were generated
134 by (i) concatenating three different sentences and then (ii) time compressing the concatenated
135 group by a factor of 3 (See Fig. 1A). Thus, each of these speeded sentence-groups had the same
136 duration as one of the original sentences. The 8 sentences used to generate the 33%-concatenated
137 (33C) condition were different than the ones used in the other conditions, and were sampled
138 independently from the Harvard sentence corpus. Ten 33C sentences were generated in this
139 manner, and each one of them was presented twice throughout the experiment, interleaved with
140 the other conditions. Compression and dilation were performed using the Overlap-Add algorithm
141 in Praat (Boersma and Weenink 2009), which preserves the spectral information of the
142 uncompressed signal.

143

144

145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167

Experimental design

Participants listened to a total of 252 sentences, divided into two blocks. Sentences were played at bedside by a laptop and speakers located in front of the patient. The experiment was controlled using Presentation® software (Neurobehavioral Systems, Inc., Berkeley, CA, www.neurobs.com). Sentences were presented in a pseudo-random order, under the constraint that the same sentence was never repeated consecutively. The experiment was self-paced: following each sentences, patients verbally rated the intelligibility of the sentence they had just heard, using a 5 point scale from 1 (“not intelligible at all”) to 5 (“fully intelligible”).

ECoG acquisition and preprocessing

Signals were recorded from 882 intracranially implanted subdural and depth electrodes (AdTech Medical Instrument Corp., WI, USA) in patients undergoing presurgical evaluation of pharmacologically intractable seizures. Electrode placement was determined solely on clinical grounds, and included grid (8×8 contacts), strip (1×4 to 1×12 contacts), and depth (1×8 contacts) electrode arrays with 10 mm inter-electrode spacing center-to-center (5 mm spacing in the depth electrodes). Neural signals were recorded on a Nicolet One EEG system, digitized at 512 Hz, and bandpass filtered between 0.5 – 250 Hz . Data were analyzed in MATLAB R2012a using custom scripts and the EEGLab toolbox (Delorme and Makeig 2004). At the preprocessing stage, each electrode was average-referenced by subtracting the mean voltage measured in all electrodes (Davidesco et al. 2013).

168 *Electrode localization*

169 Magnetic Resonance (MR) anatomical images were obtained for each patient both before and
170 after the implantation of electrodes. Electrodes were localized on the post-implant MR images
171 using intraoperative photographs, manual identification, and a custom MATLAB tool based on
172 the dimensions of the implanted electrode arrays (Yang et al. 2012). Next, the MR images were
173 nonlinearly registered to MNI space using the DARTEL algorithm in SPM (Ashburner 2007),
174 and the same transformation was used to map individual electrode coordinates into MNI space.

175

176 *Calculation of broadband power time courses*

177 Broadband power fluctuations have been shown to reflect changes in population spiking activity
178 (Crone et al. 2011; Manning et al. 2009; Nir et al. 2007; Whittingstall and Logothetis 2009). To
179 compute broadband power time courses, Morlet wavelets (standard deviation 6 cycles) with
180 center frequencies at 70, 75, 80, ... 200 Hz were convolved with the voltage time series.

181 Amplitude time series at line-noise frequencies of 120 and 180 Hz were discarded, leaving 25
182 distinct time series. Each individual amplitude time series was logarithmically transformed and
183 then converted to a z-series by subtracting its mean and dividing by its standard deviation. The
184 high frequency broadband power was then estimated as the mean of all 25 of the z-series and
185 smoothed with a hamming window of 125ms (Honey et al. 2012).

186

187 *Statistical Analysis*

188 For each one of the 882 electrodes and for each speech rate, two measures were computed:

189 1) *Neural tracking of the envelope of speech*: Neural tracking was defined as the correlation
190 between the speech envelope of a sentence and the corresponding broadband ECoG response.

191 The speech envelope was extracted for each sentence as follows: First, each sentence was filtered
192 into sixteen critical bands logarithmically spaced between 230 and 3800 Hz; second, the Hilbert
193 envelope was extracted for each band and then summed across bands. Finally, the resulting
194 envelope time-course was down-sampled to 512 Hz to match it to the ECoG signal (Doelling et
195 al. 2014). The ECoG broadband response was first shifted backwards by the response latency of
196 each electrode (as estimated from the external localizer – see below). Then, to reduce any
197 components that are not sentence-specific, the mean neural response across all sentences of a
198 given speech rate was regressed out of each trial. To account for the difference in signal length
199 across compression/dilation conditions, both the sentence envelope and the ECoG broadband
200 responses were resampled to match the original sentence duration (i.e. 33% and 66% responses
201 were up-sampled, 150% responses were down-sampled). Next, the first 300 ms and the last 300
202 ms were cropped from each trial in order to exclude onset or offset-related transients. Finally, the
203 ECoG responses were averaged across the two repetitions of each sentence and correlated with
204 the sentence envelope. Note that a correlation analysis was used, rather than a phase-locking
205 analysis, because the latter requires multiple repetitions of each sentence (Luo and Poeppel
206 2007). Repeating the same time-compressed sentence multiple times can improve its
207 intelligibility (Dupoux and Green 1997).

208 2) *Linear scaling*: This analysis was used to test the extent to which linear scaling of the neural
209 responses elicited by speeded or slowed down sentences match the neural responses to the
210 original speech rate. In this analysis, the response to the original sentence (100%) was always
211 used as a reference signal, and the neural response to each one of the other speech rates was
212 resampled to match the original sentence duration (responses to speeded speech were up-
213 sampled, responses to slowed speech were down-sampled; See Fig. 3) (Lerner et al. 2014). Then,

214 for every speech rate, the resampled neural response was correlated with the response to the
215 original sentence. Note that the linear scaling analysis is expected to provide additional
216 information, not captured by the speech tracking analysis. The speech envelope mainly reflects
217 low-level properties of the speech signal (i.e. variations in amplitude over time). High-order
218 cortical regions may no longer track the audio envelope of speech, but still be directly involved
219 in speech processing (e.g. analyzing the grammatical structure of a sentence) (Honey et al. 2012).
220 Therefore, linear scaling might be a more suitable measure to compare neural responses across
221 low- and high-level cortical areas (see Discussion).

222 In all the analyses described above, the resulting correlation values were averaged across all
223 sentences. A permutation test was used to assess the significance level of each electrode:
224 sentence labels were randomly shuffled 1000 times, such that the neural response to a given
225 sentence was correlated with the response to a different sentence. Then the empirical correlation
226 value was compared to the null distribution of correlation values in order to assess the
227 significance level of each electrode. FDR was used to correct for multiple comparisons ($q < 0.05$).

228

229 *Selection of speech-specific electrodes*

230 In the final analysis (Fig. 4), electrodes were selected based on a speech localizer task. This task
231 allowed us to contrast the mean broadband power responses to speech and to noise. In the
232 localizer task, patients viewed a still image depicting the lower part of a face, which was paired
233 with a spoken word or a noise-vocoded word (Shannon et al. 1995). There were 20-30 trials of
234 each type, and the patient was requested to press a button in response to a pre-defined target
235 word. This task was part of another experiment on audiovisual speech. Due to the limited testing
236 time available with each patient, this was the only dataset available for electrode selection.

237 However, the topography of speech-selective electrodes obtained based on this task was similar
238 to that reported by previous studies that only used auditory stimuli (Edwards et al. 2009).
239 For each trial the mean response in a time window of 50-500 ms following stimulus onset was
240 computed. Then, a t-test was used to assess whether each electrode showed a significant
241 difference in the response to speech and noise. In addition, a speech selectivity index was
242 computed for each electrode as (see Fig. 4A):

$$243 \quad \text{index} = \frac{\text{mean response to speech} - \text{mean response to noise}}{\text{mean response to speech} + \text{mean response to noise}}$$

244 A total of 40 electrodes showed a significantly stronger response to speech compared to noise
245 (False Discovery Rate (FDR) corrected, $q < 0.01$), and were thus defined as “speech-specific” and
246 used for subsequent analyses.

247 The localizer task also enabled us to extract the response latency of each electrode. The
248 Student's t-test was used to compare the broadband power response at each individual time point
249 against a pre-stimulus baseline. The response latency was defined as the time, within the time
250 series of broadband power, at which power first (i) became significantly larger than its
251 prestimulus baseline value, and (ii) remained significantly higher than baseline for at least 10
252 successive sampling points (Davidesco et al. 2013). The averaged response latency of speech-
253 selective electrodes was $150\text{ms} \pm 54\text{ms}$ (mean \pm standard deviation).

254

255

256

257

258

259

260 Results

261 *Intelligibility*

262 Patients rated the intelligibility of each sentence on a scale from 1 to 5 (where 5 is fully
263 intelligible). Intelligibility was near ceiling for the 66%, 100% and 150% rates and, as expected,
264 dropped sharply for the 33% and 33C conditions. Specifically, intelligibility significantly
265 decreased from a level of 4.88 ± 0.05 (mean rating \pm standard error of the mean) for the original
266 duration to a level of 2.78 ± 0.38 for 33% ($p=0.008$; one-sided Wilcoxon's signed-rank test), and
267 to a level of 2.19 ± 0.34 for 33C ($p=0.008$) (Fig. 1B).

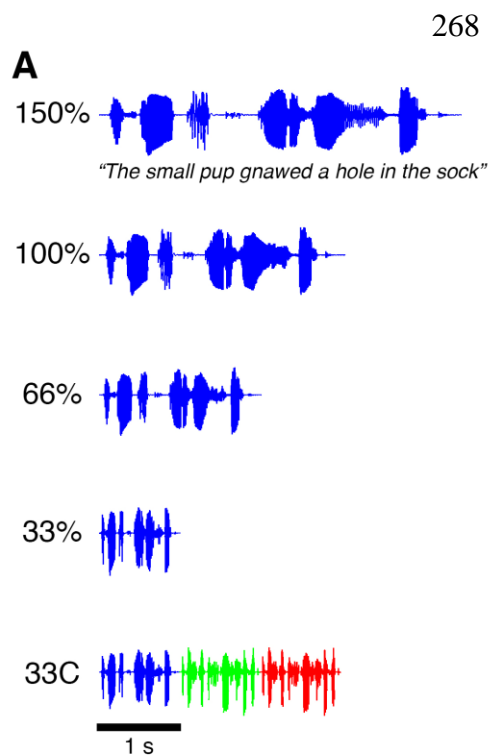
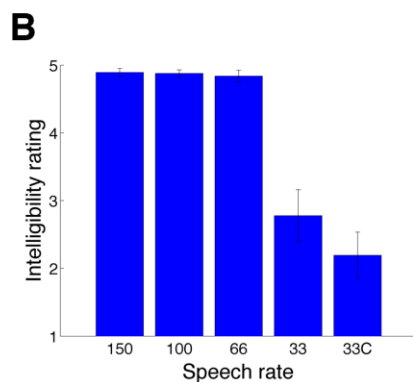


Figure 1: Experimental protocol and behavioral results

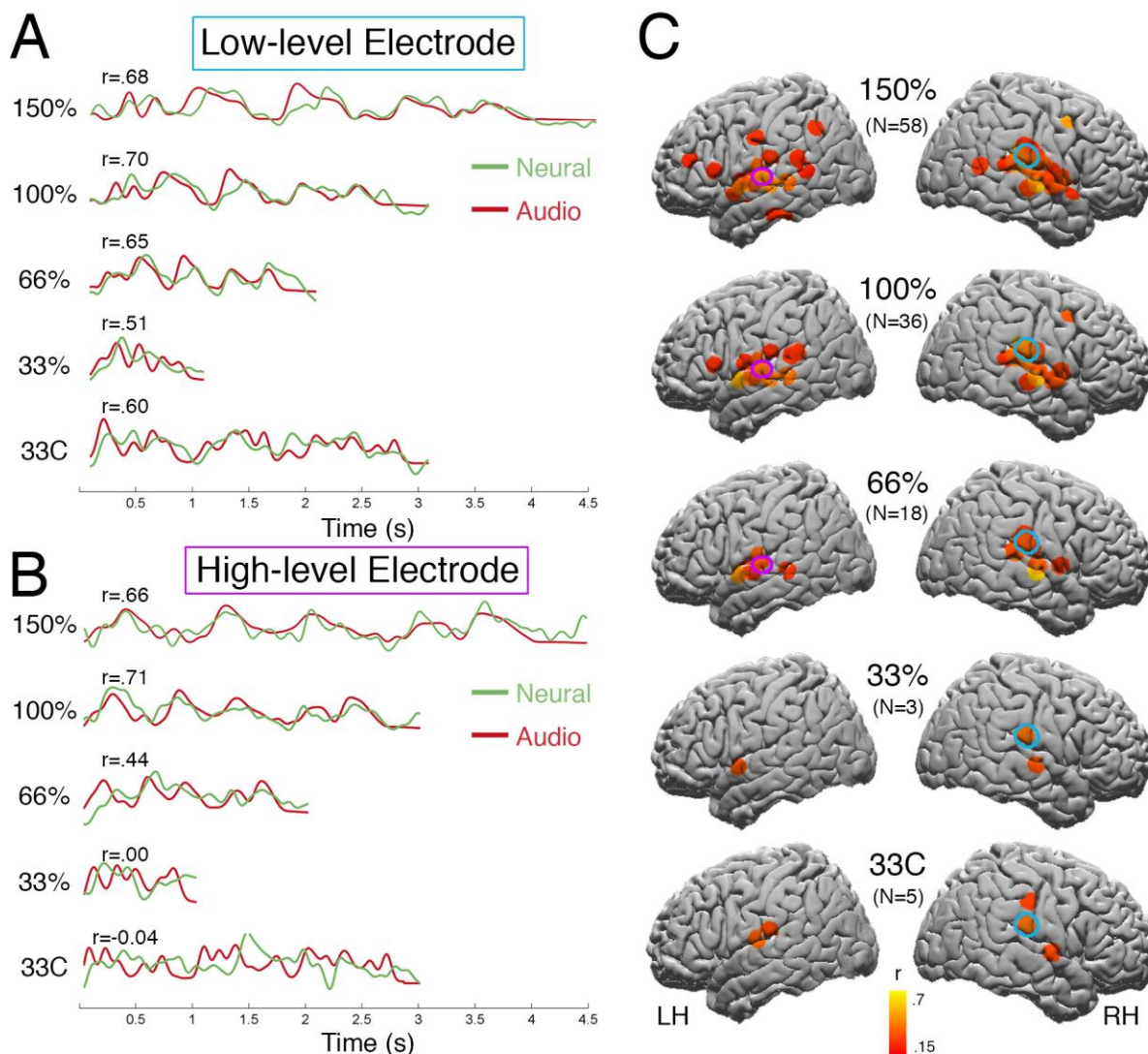
A. Experimental design: participants were presented with sentences at four different rates, with durations ranging from 33% to 150% of the original sentence duration. Participants also listened to sentences from the "33C" condition in which 3 different sentences were concatenated and then compressed by a factor of 3, to match the duration of the original sentence. **B.** Behavioral results (mean \pm SEM): intelligibility ratings were at ceiling for rates of 66% and above and dropped dramatically for 33% compression.



285 *Neural tracking*

286 We first assessed the extent to which neural responses tracked the audio envelope of each
287 sentence. Figure 2A-B depict, for two different electrodes—low-level and high-level—the
288 broadband responses (green) and audio envelope (red) for a single sentence at different speech
289 rates. The low-level auditory electrode (Fig. 2A) closely tracked the audio envelope for all
290 speech rates. In contrast, the high-level STG electrode (Fig. 2B) showed a sharp decrease in
291 envelope tracking for the 33% and 33C conditions. To investigate the spatial topography of
292 neural tracking of the speech envelope, we conducted a whole-brain analysis by assessing the
293 significance of speech tracking across all recorded electrodes using a permutation test corrected
294 for multiple comparisons (see Methods). Figure 2C shows the speech envelope tracking maps for
295 each speech rate. Even though we did not have access to neural data from Heschl’s gyrus, seven
296 electrodes, localized mainly along the lateral sulcus in the vicinity of early auditory areas,
297 displayed significant speech tracking for the most compressed speech levels (33% and 33C). For
298 intelligible speech rates (66% and above), as expected, most of the electrodes that showed
299 significant speech tracking were clustered along the STG as well as in the inferior frontal gyrus.
300 Slowed-down speech (150%) yielded the highest number of speech tracking sites, with some
301 extending to the angular gyrus and supramarginal gyrus.

302



303

304 **Figure 2: Speech envelope tracking**

305 **A.** Broadband responses of an example low-level electrode (marked in cyan in panel C) to a single sentence
306 presented at different rates (green) as well as the audio envelope of that sentence (red). Correlation values
307 correspond to the single sentence that is depicted in panel A.

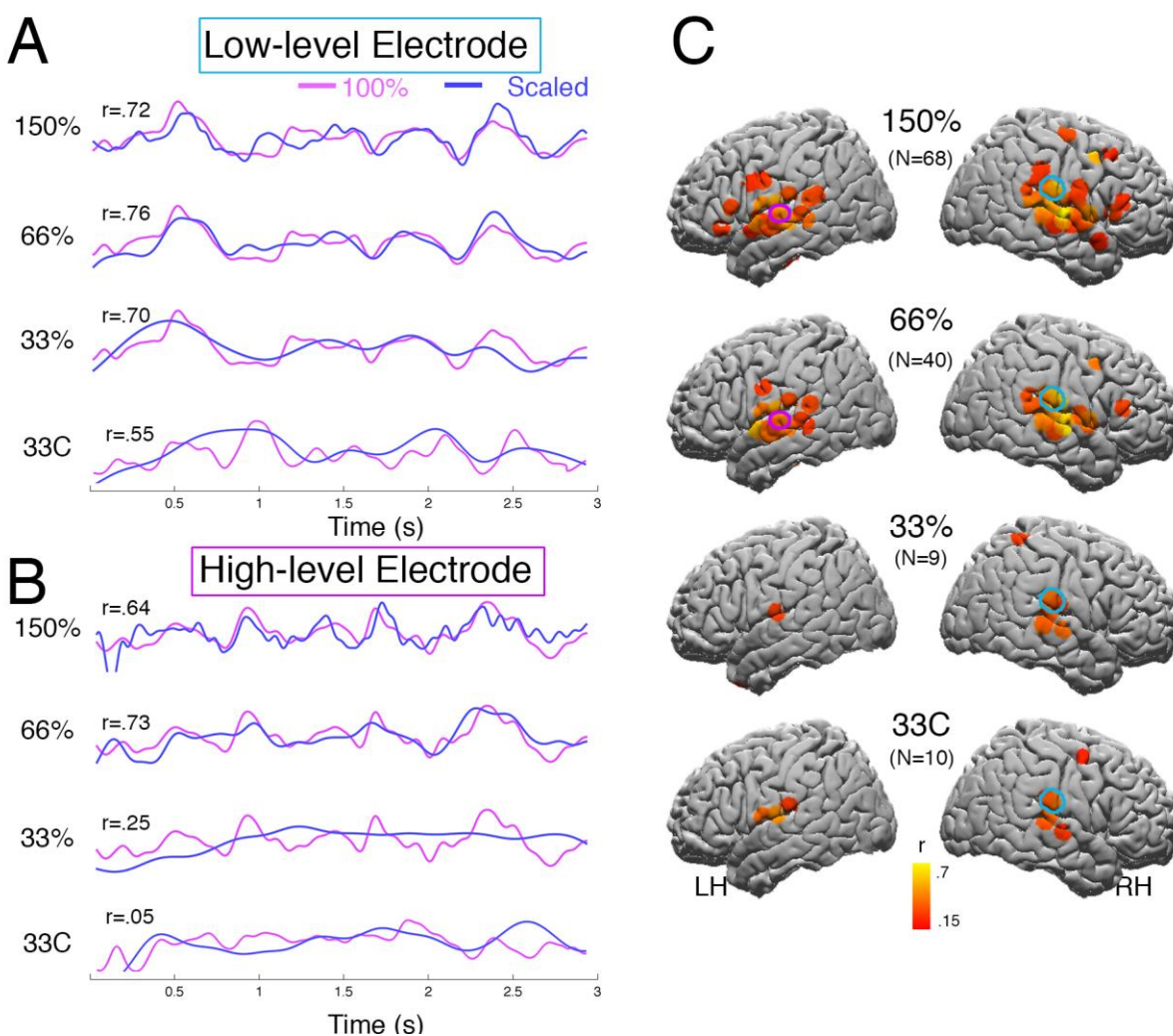
308 **B.** Same as (A) for an example high-level electrode (marked in purple in panel C).

309 **C.** Spatial distribution of all the electrodes that showed significant speech tracking ($p < 0.05$, FDR corrected) at each
310 one of the speech rates. Electrodes are color-coded based on the correlation of the broadband response with speech
311 audio envelope. Significance was assessed using a permutation test.

312

313

314 Could these results be driven by the difference in signal duration across conditions? To address
315 that question, we compared the 33C and 100% conditions. These two conditions differ only by
316 speech rate, not by signal length (see Methods). The 100% condition was characterized by wide-
317 spread speech tracking along the anterior and posterior STG as well as the inferior frontal gyrus,
318 whereas in the 33C condition, only 5 electrodes in the vicinity of early auditory areas exhibited
319 significant speech tracking. Moreover, prior to measuring speech tracking, both the sentence
320 envelope and the ECoG broadband responses for the compressed/dilated conditions were
321 resampled to match the original sentence duration, therefore the number of time points in the
322 neural response was held constant across speech rates (see Methods).



323

324 **Figure 3: Temporal scaling**

325 **A.** Broadband responses of an example low-level electrode (same electrode as in Fig. 2A). The neural responses to
326 speeded (slowed) sentences (blue) were up-sampled (down-sampled) to match the length of the neural response to
327 the original sentence (pink).

328 **B.** Broadband responses of an example high-level electrode (same electrode as in Fig. 2B).

329 **C.** Spatial distribution of electrodes showing significant temporal scaling ($p < 0.05$; FDR corrected) for each speech
330 rate.

331

332 *Linear scaling*

333 Based on the work of Lerner et al. (2014), we also examined the linear scaling of the neural
334 responses. Here, we measured the extent to which the responses for the speeded (slowed)
335 sentences match the original (100%) by up-sampling (down-sampling) the neural responses (see
336 Methods) (Lerner et al. 2014).

337 Figure 3A-B depict the scaled neural responses (blue) and the response to the original sentence
338 that served as a reference (pink). For example, in the case of 150% (slowed) speech, the neural
339 response was compressed (i.e. down-sampled) to match the 100% response, and then the two
340 responses were correlated.

341 Similarly to the speech tracking analysis, we identified two representative electrodes. For low-
342 level auditory electrodes (Fig. 3A), significant response scaling was observed across all speech
343 rates, even outside the intelligibility range. In contrast, for high-order STG electrodes (Fig. 3B),
344 temporal scaling was observed only for intelligible speech (66% and 150%) and not for
345 unintelligible speech (33% and 33C). This step-like transition in temporal scaling from
346 intelligible to non-intelligible speech was also evident in a whole-brain analysis. Significant
347 temporal scaling along STG, Inferior Frontal Gyrus (IFG) and supramarginal gyrus was observed
348 for intelligible speech (66% and 150%). In contrast, scaling of neural responses to unintelligible

349 speech was mainly confined to STG sites in close proximity to early auditory cortex.
350 Finally, to further explore how speech rate affects neural processing in language related areas,
351 we focused our analysis on 40 electrodes, which exhibited increased neural response to speech
352 relative to non-speech stimuli, defined using an independent localizer task (see Methods). These
353 electrodes were mainly clustered along the right and left STG, with the exception of 4 electrodes
354 that were distributed over the IFG and motor cortex (Fig. 4A). Across the 40 speech-specific
355 electrodes, the correlation between broadband responses and audio envelope (Fig. 4B) decreased
356 monotonically as speech rate increased. Unlike speech envelope tracking, temporal scaling
357 values dropped sharply, in a step-like function, for non-intelligible speech, in accordance with
358 intelligibility ratings (compare Fig. 4C to Fig. 1B). To directly assess the transition from
359 intelligible to non-intelligible speech across metrics, we conducted a two-way repeated measures
360 ANOVA with speech ratio (66% vs. 33%) and metric (envelope tracking vs. temporal scaling) as
361 factors. There was a highly significant interaction between these two factors ($F(1,39)=46.80$,
362 $p<10^{-9}$), indicating that temporal scaling dropped significantly in the transition between 66% and
363 33%, whereas speech tracking did not. This suggests that temporal scaling might be a more
364 sensitive measure of speech processing that is more closely related to the observed behavioral
365 effect.

366

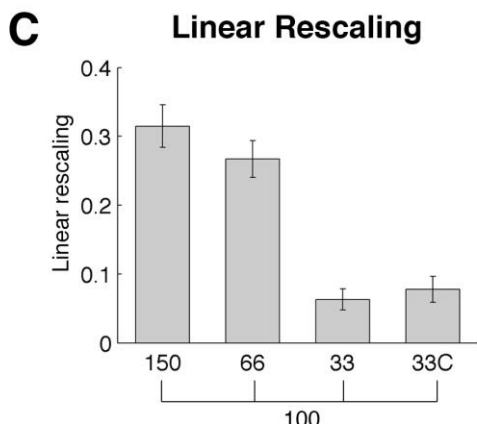
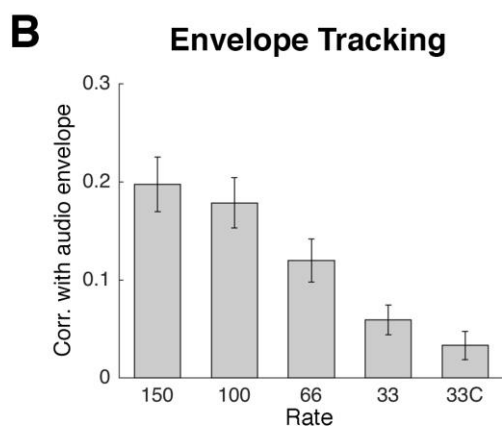
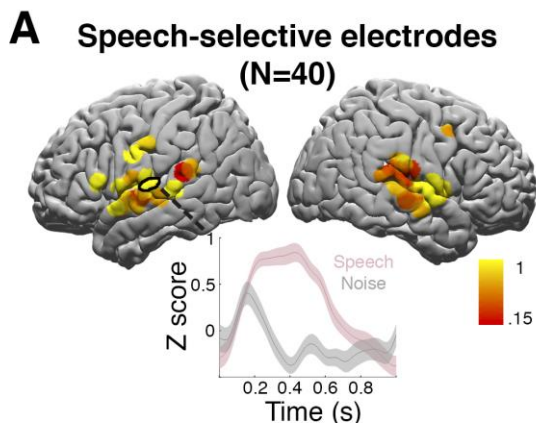


Figure 4: Response profile of speech-specific electrodes

A. Speech specificity map: electrodes that showed a significantly stronger response to individual words compared to noise-vocoded words in an independent localizer task ($p < 0.05$, FDR corrected). Electrodes are color coded according to speech specificity (0 = non selective, 1 = highly selective). Inset shows the broadband (70-200 Hz) responses to speech (red) and noise (black) of an example speech-selective electrode.

B. Speech tracking correlation values (mean \pm SEM) across 40 speech-specific electrodes as a function of speech rate.

C. Same as (B) for linear scaling values.

386 **Discussion**

387 The human auditory system can comprehend spoken language with remarkable tolerance to
388 speech rate. Such tolerance, however, is limited. In particular, it has long been known that
389 artificially time-compressing speech to a level beyond what is normally encountered in everyday
390 listening (e.g., beyond compression by 3) hinders intelligibility at the word level and
391 comprehension at the sentence level (Dupoux and Green 1997; Foulke and Sticht 1969; Garvey
392 1953; Ghitza 2014). Here, we recorded ECoG responses to sentences presented at speeded rates
393 (33%, 66%), at the original rate (100%) and at a slowed rate (150% of the original sentence
394 duration). Behaviorally, patients reported much lower understanding of highly compressed
395 speech (33%; Fig. 1B). At the neural level, we observed two distinct response profiles: 1) A
396 ‘low-level’ profile in electrodes along the STG, adjacent A1+, in which envelope tracking and
397 linear scaling were observed across all speech rates; 2) A ‘high-level’ profile in electrodes
398 further along the cortical hierarchy in anterior STG, posterior STG, IFG and Supramarginal
399 gyrus, in which we observed significant speech tracking and linear scaling only within the
400 intelligibility range (Figs. 2 and 3). Most of the electrodes in the current study exhibited the
401 ‘high-level’ response profile, potentially as a result of the limited electrode coverage (i.e. we did
402 not have access to recordings directly from Heschl's Gyrus).

403 Our results help reconcile seemingly contradictory findings in the literature. Nourski et al.
404 (2009), using intracranial recordings, demonstrated that Heschl's gyrus (primary auditory cortex)
405 can track the speech envelope well outside the intelligibility range. On the other hand, Ahissar et
406 al. (2001), using MEG, reported that time compression of speech beyond the intelligibility limit
407 is associated with a sharp decrease in speech envelope tracking. Our results suggest that these
408 previous findings might correspond to distinct processing stages along the cortical processing

409 hierarchy. Even though we did not have access to recordings in Heschl's gyrus within the Sylvian
410 fissure, adjacent areas along the STG provided similar findings to these reported by Nourski et
411 al. (2009). In contrast, the response in higher-order linguistic and extra-linguistic areas along the
412 STG, IFG and supramarginal gyrus exhibited similar response profile to that reported by Ahissar
413 et al. (2001) and Lerner et al. (2014). These results are consistent with a recent intracranial EEG
414 study, which demonstrated a hierarchical organization of sound processing from the primary
415 auditory cortex, where activity closely reflects the acoustic features of the stimulus, through the
416 STG, where activity reflects both acoustic features and task demands, to the prefrontal cortex,
417 which is mainly modulated by task requirements and behavioral performance (Nourski 2017).
418 A related question that has received attention in the literature is where along the cortical
419 hierarchy is the bottleneck in the processing of time-compressed speech. Our results, in
420 accordance with Nourski et al. (2009), demonstrate that early auditory cortex can track speech
421 outside of the intelligibility range, and it is therefore not to be considered the bottleneck. In
422 accordance with our hypothesis, the first sites to track the speech envelope and to exhibit neural
423 scaling only for intelligible speech were located along the STG. Interestingly, these areas seem to
424 have an intermediate processing timescale in the order of few hundreds of milliseconds, which
425 corresponds to the formation of syllables and the integration of syllables into words (Hasson et
426 al. 2008; Honey et al. 2012; Lerner et al. 2011). Given that information flows upstream along the
427 timescale hierarchy (from early auditory cortex to linguistic and extra-linguistic regions), it is
428 reasonable to hypothesize that the bottleneck lies in areas with relatively short TRW that decode
429 syllables and integrate syllables into words.

430 Once acoustic information is integrated into words, it can be transmitted up the timescale
431 hierarchy to areas with longer TRW needed for the integration of words into sentences and

432 sentences into paragraphs. Indeed, we observed that the neural activity in high-level linguistic
433 areas in the IFG and supramarginal gyrus tracked the speech envelope only for intelligible
434 speech. These findings are congruent with an fMRI study that demonstrated that the inferior
435 frontal gyrus and the superior temporal sulcus show an invariant response to moderate
436 compression rates followed by a sharp decline in activation for non-intelligible compressed
437 speech (Vagharchakian et al., 2012). The current study extends these findings by demonstrating
438 that millisecond-by-millisecond STG responses track the speech envelope and linearly scale with
439 speech rate, as long as speech remains intelligible. Furthermore, our finding that language areas
440 scale their dynamics in response to speech rate suggests that temporal integration windows
441 should also be assessed using relative information-based units (e.g. the number of syllables)
442 rather than merely in absolute temporal units (e.g. milliseconds), which vary across compression
443 rates (Lerner et al. 2014).

444 It is worth noting that due to the limited testing time available with neurosurgical patients, this
445 study only examined four speech rates. In future studies, it would be informative to sample the
446 intelligibility spectrum more densely in order to map the transition from intelligible to non-
447 intelligible speech.

448 Why do areas with an intermediate TRW, which are presumably involved in syllable formation
449 and the integration of syllables into words, fail to track speech compressed by a factor of 3 or
450 more? A potential explanation is provided by TEMPO (Ghitza 2011), a model that epitomizes
451 recently proposed oscillation-based models of speech perception (Ahissar and Ahissar 2005;
452 Ding and Simon 2009; Ghitza and Greenberg 2009; Giraud and Poeppel 2012; Hyafil et al. 2015;
453 Lakatos et al. 2005; Peelle and Davis 2012; Poeppel 2003). TEMPO postulates a cortical
454 computation principle by which decoding is performed within a hierarchical time-varying

455 window structure, synchronized with the input on multiple time scales. The windows are
456 generated by a segmentation process, implemented by a cascade of oscillators, governed by the
457 theta oscillator, which provides syllabic segmentation. These oscillators operate within a
458 constrained range of frequencies (the biophysical frequency range of theta). Critically,
459 intelligibility remains high as long as theta is in sync with the input (as is the case for moderate
460 speech speeds) and it sharply deteriorates once theta is out of sync (when the input syllabic rate
461 is outside the theta frequency range). The notion that cortical oscillations are closely related to
462 speech uptake capacity has received support from several recent studies (Borges et al. 2018;
463 Pefkou et al. 2017). The findings of the current study suggest that the neuronal circuitry of the
464 theta oscillator might be located at the STG level.

465 Finally, it is worth noting the difference between the insights provided by the neural tracking and
466 the linear scaling measures. While neural tracking measures how well the neural response
467 matches the acoustic envelope, linear scaling captures how consistent the neural response is
468 across speech rates. Neural tracking is mostly sensitive to low-level properties of the speech
469 signal (i.e. variations in amplitude across time). Whereas low-level regions (e.g. A1+) are
470 expected to closely track the speech envelope, this might not be the case with high-order cortical
471 regions. Indeed, Honey et al. (2012) demonstrated that regions with long TRW (e.g. medial
472 frontal gyrus) no longer track the audio envelope, and yet respond very reliably to audiovisual
473 stimuli. In the current study, more electrodes showed significant linear scaling compared to
474 speech tracking (e.g. in the 66% rate: 18 electrodes compared 40 electrodes, respectively).

475 Moreover, note that the degree of linear scaling dropped sharply with speech rate, in a better
476 correspondence to the behavioral data, compared to a shallower drop in neural tracking (Fig. 4).

477 This finding extends Lerner et al. (2014), who observed a linear scaling effect in temporally

478 sluggish fMRI measurements. Here, we show that millisecond-by-millisecond neural responses
479 recorded directly from high-level auditory regions can linearly scale with speech rate, as long as
480 speech remains intelligible. Our findings suggest that neural tracking at secondary auditory areas
481 in the STG, and beyond, is a prerequisite for intelligibility. As long as envelope tracking is
482 maintained, the linearly scaled neural responses are remarkably stable as speech rate varies, and
483 speech is intelligible.

484

485 **Acknowledgements**

486 This work was supported by US National Institutes of Health grant R01 MH094480 (U.H.,
487 C.J.H.). We would like to acknowledge the contribution of the patients who participated in this
488 study.

489

490 **Disclosures**

491 The authors declare no competing financial interests

492

493

494

495

496

497

498

499

500 **References**

- 501 **Ahissar E, and Ahissar M.** Processing of the temporal envelope of speech. In: *The auditory*
502 *cortex: A synthesis of human and animal research*, edited by Peter Heil HS, Eike Budinger,
503 Reinhard Konig. London: Lawrence Erlbaum Associates, Inc, 2005, p. 295-313.
- 504 **Ahissar E, Nagarajan S, Ahissar M, Protopapas A, Mahncke H, and Merzenich MM.**
505 Speech comprehension is correlated with temporal response patterns recorded from auditory
506 cortex. *Proceedings of the National Academy of Sciences* 98: 13367-13372, 2001.
- 507 **Ashburner J.** A fast diffeomorphic image registration algorithm. *Neuroimage* 38: 95-113, 2007.
- 508 **Boersma P, and Weenink D.** Praat: doing phonetics by computer (Version 5.1. 05)[Computer
509 program]. 2009.
- 510 **Borges AFT, Giraud A-L, Mansvelder HD, and Linkenkaer-Hansen K.** Scale-free amplitude
511 modulation of neuronal oscillations tracks comprehension of accelerated speech. *Journal of*
512 *Neuroscience* 38: 710-722, 2018.
- 513 **Crone NE, Korzeniewska A, and Franaszczuk PJ.** Cortical gamma responses: searching high
514 and low. *International Journal of Psychophysiology* 79: 9-15, 2011.
- 515 **Davidesco I, Zion-Golumbic E, Bickel S, Harel M, Groppe DM, Keller CJ, Schevon CA,**
516 **McKhann GM, Goodman RR, and Goelman G.** Exemplar selectivity reflects perceptual
517 similarities in the human fusiform cortex. *Cerebral cortex* 24: 1879-1893, 2013.
- 518 **Delorme A, and Makeig S.** EEGLAB: an open source toolbox for analysis of single-trial EEG
519 dynamics including independent component analysis. *Journal of neuroscience methods* 134: 9-
520 21, 2004.
- 521 **Ding N, and Simon JZ.** Neural representations of complex temporal modulations in the human
522 auditory cortex. *Journal of neurophysiology* 102: 2731-2743, 2009.

- 523 **Doelling KB, Arnal LH, Ghitza O, and Poeppel D.** Acoustic landmarks drive delta–theta
524 oscillations to enable speech comprehension by facilitating perceptual parsing. *Neuroimage* 85:
525 761-768, 2014.
- 526 **Dupoux E, and Green K.** Perceptual adjustment to highly compressed speech: effects of talker
527 and rate changes. *Journal of Experimental Psychology: Human perception and performance* 23:
528 914, 1997.
- 529 **Edwards E, Soltani M, Kim W, Dalal SS, Nagarajan SS, Berger MS, and Knight RT.**
530 Comparison of time–frequency responses and the event-related potential to auditory speech
531 stimuli in human cortex. *Journal of neurophysiology* 102: 377-386, 2009.
- 532 **Foulke E, and Sticht TG.** Review of research on the intelligibility and comprehension of
533 accelerated speech. *Psychological bulletin* 72: 50, 1969.
- 534 **Garvey WD.** The intelligibility of speeded speech. *Journal of Experimental Psychology* 45: 102,
535 1953.
- 536 **Ghitza O.** Behavioral evidence for the role of cortical θ oscillations in determining auditory
537 channel capacity for speech. *Frontiers in psychology* 5: 652, 2014.
- 538 **Ghitza O.** Linking speech perception and neurophysiology: speech decoding guided by cascaded
539 oscillators locked to the input rhythm. *Frontiers in psychology* 2: 130, 2011.
- 540 **Ghitza O, and Greenberg S.** On the possible role of brain rhythms in speech perception:
541 intelligibility of time-compressed speech with periodic and aperiodic insertions of silence.
542 *Phonetica* 66: 113-126, 2009.
- 543 **Giraud A-L, and Poeppel D.** Cortical oscillations and speech processing: emerging
544 computational principles and operations. *Nature neuroscience* 15: 511, 2012.

545 **Hasson U, Yang E, Vallines I, Heeger DJ, and Rubin N.** A hierarchy of temporal receptive
546 windows in human cortex. *Journal of Neuroscience* 28: 2539-2550, 2008.

547 **Honey CJ, Thesen T, Donner TH, Silbert LJ, Carlson CE, Devinsky O, Doyle WK, Rubin**
548 **N, Heeger DJ, and Hasson U.** Slow cortical dynamics and the accumulation of information over
549 long timescales. *Neuron* 76: 423-434, 2012.

550 **Hyafil A, Fontolan L, Kabdebon C, Gutkin B, and Giraud A-L.** Speech encoding by coupled
551 cortical theta and gamma oscillations. *Elife* 4: 2015.

552 **IEEE.** IEEE Recommended Practice for Speech Quality Measurements. In: *IEEE1969*, p. 1-24.

553 **Lakatos P, Shah AS, Knuth KH, Ulbert I, Karmos G, and Schroeder CE.** An oscillatory
554 hierarchy controlling neuronal excitability and stimulus processing in the auditory cortex.
555 *Journal of neurophysiology* 94: 1904-1911, 2005.

556 **Lerner Y, Honey CJ, Katkov M, and Hasson U.** Temporal scaling of neural responses to
557 compressed and dilated natural speech. *Journal of neurophysiology* 111: 2433-2444, 2014.

558 **Lerner Y, Honey CJ, Silbert LJ, and Hasson U.** Topographic mapping of a hierarchy of
559 temporal receptive windows using a narrated story. *Journal of Neuroscience* 31: 2906-2915,
560 2011.

561 **Luo H, and Poeppel D.** Phase patterns of neuronal responses reliably discriminate speech in
562 human auditory cortex. *Neuron* 54: 1001-1010, 2007.

563 **Manning JR, Jacobs J, Fried I, and Kahana MJ.** Broadband shifts in local field potential
564 power spectra are correlated with single-neuron spiking in humans. *Journal of Neuroscience* 29:
565 13613-13620, 2009.

- 566 **Mehler J, Sebastian N, Altmann G, Dupoux E, Christophe A, and Pallier C.** Understanding
567 compressed sentences: The role of rhythm and meaning. *Annals of the New York Academy of*
568 *Sciences* 682: 272-282, 1993.
- 569 **Nir Y, Fisch L, Mukamel R, Gelbard-Sagiv H, Arieli A, Fried I, and Malach R.** Coupling
570 between neuronal firing rate, gamma LFP, and BOLD fMRI is related to interneuronal
571 correlations. *Current biology* 17: 1275-1285, 2007.
- 572 **Nourski KV.** Auditory processing in the human cortex: An intracranial electrophysiology
573 perspective. *Laryngoscope investigative otolaryngology* 2: 147-156, 2017.
- 574 **Nourski KV, Reale RA, Oya H, Kawasaki H, Kovach CK, Chen H, Howard MA, and**
575 **Brugge JF.** Temporal envelope of time-compressed speech represented in the human auditory
576 cortex. *Journal of Neuroscience* 29: 15564-15574, 2009.
- 577 **Pallier C, Sebastian-Gallés N, Dupoux E, Christophe A, and Mehler J.** Perceptual
578 adjustment to time-compressed speech: A cross-linguistic study. *Memory & cognition* 26: 844-
579 851, 1998.
- 580 **Peelle JE, and Davis MH.** Neural oscillations carry speech rhythm through to comprehension.
581 *Frontiers in psychology* 3: 320, 2012.
- 582 **Pefkou M, Arnal LH, Fontolan L, and Giraud A-L.** θ -Band and β -Band Neural Activity
583 Reflects Independent Syllable Tracking and Comprehension of Time-Compressed Speech.
584 *Journal of Neuroscience* 37: 7930-7938, 2017.
- 585 **Poeppel D.** The analysis of speech in different temporal integration windows: cerebral
586 lateralization as ‘asymmetric sampling in time’. *Speech communication* 41: 245-255, 2003.
- 587 **Sebastián-Gallés N, Dupoux E, Costa A, and Mehler J.** Adaptation to time-compressed
588 speech: Phonological determinants. *Perception & psychophysics* 62: 834-842, 2000.

589 **Shannon RV, Zeng F-G, Kamath V, Wygonski J, and Ekelid M.** Speech recognition with
590 primarily temporal cues. *Science* 270: 303-304, 1995.

591 **Vagharchakian L, Dehaene-Lambertz G, Pallier C, and Dehaene S.** A temporal bottleneck in
592 the language comprehension network. *Journal of Neuroscience* 32: 9089-9102, 2012.

593 **Whittingstall K, and Logothetis NK.** Frequency-band coupling in surface EEG reflects spiking
594 activity in monkey visual cortex. *Neuron* 64: 281-289, 2009.

595 **Yang AI, Wang X, Doyle WK, Halgren E, Carlson C, Belcher TL, Cash SS, Devinsky O,**
596 **and Thesen T.** Localization of dense intracranial electrode arrays using magnetic resonance
597 imaging. *Neuroimage* 63: 157-165, 2012.

598