

1 **Title: Pre- and post-target cortical processes predict speech-in-noise performance**

2 Abbreviated title: Neural correlates of speech-in-noise performance

3 **Subong Kim,^{1, #} Adam T. Schwalje,^{2, #} Andrew S. Liu,² Phillip E. Gander,³ Bob McMurray,**
4 **^{1,2,4} Timothy D. Griffiths,⁵ and Inyong Choi^{1, 2, *}**

5 ¹ Department of Communication Sciences and Disorders, University of Iowa, Iowa City, IA
6 52242, USA

7 ² Department of Otolaryngology – Head and Neck Surgery, University of Iowa Hospitals and
8 Clinics, Iowa City, IA 52242, USA

9 ³ Department of Neurosurgery, University of Iowa Hospitals and Clinics, Iowa City, IA 52242,
10 USA

11 ⁴ Department of Psychological and Brain Sciences, University of Iowa, Iowa City, IA 52242, USA

12 ⁵ Biosciences Institute, Newcastle University, Newcastle upon Tyne, NE1 7RU, UK

13 # SK and ATS contributed equally to this work.

14 * Corresponding author at: Department of Communication Sciences and Disorders, University of
15 Iowa, 250 Hawkins Dr., Iowa City, IA 52242, USA. Email address: inyong-choi@uiowa.edu

16

17 **Author contributions:**

18 S.K. and A.T.S. contributed equally to this work. I.C. designed the experiments. S.K. ran the
19 experiments. S.K., A.T.S., A.S.L., P.E.G., B.M., T.D.G., and I.C. analyzed and interpreted data.

20 S.K., A.T.S., and I.C. prepared the manuscript, with revisions and suggestions from A.S.L.,
21 P.E.G., B.M., and T.D.G..

22

23

24 **Abstract**

25 Understanding speech in noise (SiN) is a complex task that recruits multiple cortical
26 subsystems. There is a variance in individuals' ability to understand SiN that cannot be
27 explained by simple hearing profiles, which suggests that central factors may underlie the
28 variance in SiN ability. Here, we elucidated a few cortical functions involved during a SiN task
29 and their contributions to individual variance using both within- and across-subject approaches.
30 Through our within-subject analysis of source-localized electroencephalography, we
31 investigated how acoustic signal-to-noise ratio (SNR) alters cortical evoked responses to a
32 target word across the speech recognition areas, finding stronger responses in left
33 supramarginal gyrus (SMG, BA40 the *dorsal lexicon* area) with quieter noise. Through an
34 individual differences approach, we found that listeners show different neural sensitivity to the
35 background noise and target speech, reflected in the amplitude ratio of earlier auditory-cortical
36 responses to speech and noise, named as an *internal SNR*. Listeners with better *internal SNR*
37 showed better SiN performance. Further, we found that the post-speech time SMG activity
38 explains a further amount of variance in SiN performance that is not accounted for by *internal*
39 *SNR*. This result demonstrates that at least two cortical processes contribute to SiN
40 performance independently: pre-target time processing to attenuate neural representation of
41 background noise and post-target time processing to extract information from speech sounds.

42

43 **Keywords**

44 Speech-in-noise, speech unmasking, speech recognition, individual differences,
45 electroencephalography, supramarginal gyrus

46

47 1. Introduction

48 Understanding speech in noise (SiN) is essential for communication in social settings.
49 Young normal-hearing listeners are remarkably adept at this. Even in challenging SiN conditions
50 where the speech and noise have the same intensity (i.e., 0 dB signal-to-noise ratio: SNR) and
51 overlapped frequency components, they often recognize nearly 90% of sentences correctly
52 (Ohlenforst et al., 2017). This suggests a surprising capacity of the auditory system to cope with
53 noise. However, the ability to understand SiN degrades severely with increased background
54 noise level (Ohlenforst et al., 2017), hearing loss (Harris and Swenson, 1990), and/or aging
55 (Nabelek, 1988).

56 Recent studies show that normal hearing listeners show large individual differences in
57 SiN performance (Lieberman et al., 2016). The premise of this study is that by linking this
58 variable ability for SiN perception to variation in cortical activity, we may be able to understand
59 the neural mechanisms by which humans accomplish this ability, and this may shape our
60 understanding of how best to remediate hearing loss.

61 Two broad neural mechanisms might give rise to better or worse SiN performance. First,
62 listeners may vary in neural processing that separates the target auditory object from the
63 mixture of sounds (i.e., similar to *external* selection processes in (Strauss and Francis, 2017).
64 Auditory scene analysis (Bregman, 1999) processes, when occurring in parallel with auditory
65 selective attention (Shinn-Cunningham, 2020), can inhibit the neural representation of
66 competing sounds and enhance the neural response to attended input. This process has been
67 conceptualized as a form of sensory gain control (Hillyard et al., 1998) or neural filtering
68 (Obleser and Erb, 2020). Its effectiveness is often quantified as attentional modulation index
69 (AMI), the amplitude ratio of evoked responses to background noise and target (Dai and Shinn-
70 Cunningham, 2016; O'Sullivan et al., 2019), or the degree of neural phase-locking to the
71 attended speech (Etard and Reichenbach, 2019; Mesgarani and Chang, 2012; Viswanathan et
72 al., 2019). A successful sensory gain control, indicated by a positive AMI, during a SiN task will

73 unmask the target speech from maskers, which will enhance the effective signal-to-noise ratio
74 (SNR) in the central auditory pathway (e.g., the primary and secondary auditory cortices in the
75 superior temporal plane: STP and the posterior superior temporal gyrus: STG).

76 Second, listeners might vary in neural processes for the prompt extraction of information
77 from a speech signal (i.e., similar to *internal* selection processes in (Strauss and Francis, 2017)).
78 An inherent challenge in recognizing speech (e.g., a spoken word) is the mapping between the
79 incoming speech cues and higher level units like words and meaning while speech unfolds
80 rapidly over time [for review, see Weber and Scharenborg (2012), Dahan and Magnuson
81 (2006), and Davis (2016) with references therein]. In quiet listening conditions, average young
82 normal-hearing listeners activate a range of lexical candidates immediately at the onset of the
83 auditory stimulus (i.e., shown by works using eye-movements in the visual world paradigm:
84 (Allopenna et al., 1998; Dahan and Gareth Gaskell, 2007; Magnuson et al., 2007). For example,
85 after hearing the /ba/ at the onset of *bakery*, listeners will immediately consider a range of words
86 like *bacon*, *bathe*, or *base*, at both phonological and semantic levels. However, such rapid
87 lexical processing develops slowly in children (Rigler et al., 2015); continuous differences in
88 lexical processing are linked to differences in language ability (McMurray et al., 2010); and they
89 differ in listeners with hearing loss or deteriorated acoustic-cue encoding (McMurray et al.,
90 2019). These facts suggest that there can be individual differences in the prompt lexical
91 processing even for *clean* speech signals.

92 It is as yet unclear the degree to which variation in SiN performance is related to
93 variation in both processes (particularly in combination).

94 ***Assessing Individual Differences in Speech Unmasking.*** Individual differences in the
95 speech unmasking pathway may arise from 1) the fidelity of encoding supra-threshold acoustic
96 features and 2) cognitive control of the domain-general attentional network. Auditory scene
97 analysis relies on the supra-threshold acoustic features that provide binding cues for auditory
98 grouping (Darwin, 1997). These include the spectra (Lee et al., 2013), location (Frey et al.,

99 2014; Goldberg et al., 2014), temporal coherence (Moore, 1990; Shamma et al., 2013; Teki et
100 al., 2011), rhythm (Calderone et al., 2014; Golumbic et al., 2013; Herrmann et al., 2016; Obleser
101 and Kayser, 2019), and timing (Lange, 2009) of the figure and ground. The fidelity of encoding
102 such supra-threshold acoustic features may affect the separation of target speech from
103 background noise. Supporting this idea, previous studies have correlated the fidelity of supra-
104 threshold acoustic cue coding to SiN understanding (Anderson and Kraus, 2010; Anderson et
105 al., 2013; Holmes and Griffiths, 2019; Hornickel et al., 2009; Liberman et al., 2016; Parbery-
106 Clark et al., 2009; Song et al., 2011).

107 Individual differences also exist in how strongly selective attention modulates cortical
108 evoked responses to sounds (Choi et al., 2014). This suggests there may be a correlation
109 between top-down selective attention efficacy and SiN performance. Indeed, (Strait and Kraus,
110 2011) reported that reaction time during a selective attention task predicts SiN performance.
111 Similarly, studies suggest that poor cognitive control of executive attentional network predicts
112 auditory selective attention performance (Bressler et al., 2017; Dai et al., 2018), though this has
113 not been extended to SiN.

114 We can obtain a measure of the overall function of these bottom-up and top-down neural
115 processing for speech unmasking by quantifying the amplitude ratio of early cortical auditory
116 evoked responses to noise and target speech (similarly to the AMI concept). Here we use the
117 N1/P2 event-related potential (ERP) components which occur with 100-300 ms latency.
118 Previous studies showed that such ERP components are strongly modulated by selective
119 attention but only when auditory objects are successfully segregated (Choi et al., 2013; Choi et
120 al., 2014; Kong et al., 2015). Since those early cortical ERP components originate from multiple
121 regions across Heschl's gyrus (i.e., the primary auditory cortex) and its surrounding areas (e.g.,
122 posterior superior temporal gyrus) (Céronien et al., 1998), an efficient and collective way of
123 indexing the neural efficiency in speech unmasking is using a scalp electroencephalographical
124 (EEG) potential at the vertex [e.g., "Cz" of the international 10-10 system for EEG electrode

125 montage: Koessler et al. (2009) within a limited time-window (e.g., 100-300 ms range after the
126 stimulus onset].

127 ***Assessing Individual Differences in Mapping Speech to Words and Meaning.*** To

128 assess individual differences related to the second neural mechanism – the downstream speech
129 information processing – we must assess a larger range of cortical regions above the auditory
130 brainstem and cortex. Current models of speech processing suggest two distinct cortical
131 networks (i.e., dorsal and ventral stream) that are used in parallel (Gow, 2012; Hickok and
132 Poeppel, 2007; Myers et al., 2009; Scott and Johnsrude, 2003). The ventral stream pathway
133 including anterior superior temporal and middle temporal gyri (STG/MTG) integrates speech-
134 acoustic and semantic information progressively over time for the sound-to-meaning mapping
135 (Davis and Johnsrude, 2003). The dorsal stream pathway comprising supramarginal gyrus
136 (SMG, also known as tempo-parietal junction or TPJ) and pre- / post-central gyri mediates the
137 mapping between sound and articulation (Rauschecker and Scott, 2009), while inferior frontal
138 gyrus (IFG) interacts with both pathways for lexical decision-making processes (Gow, 2012).

139 Studies have compared cortical responses in these areas to spoken words against
140 acoustically-matching non-word sounds. These highlight the SMG/TPJ, MTG, and IFG in the
141 left hemisphere as three regions that tend to exhibit more activity for words than pseudo-words
142 (Davis and Gaskell, 2009; Taylor et al., 2013). The dual lexicon model suggested by Gow
143 (2012) confirms the importance of those three regions by referring to left SMG and MTG as
144 dorsal and ventral lexicons that communicate with left IFG for lexical decision making. While
145 both SMG and MTG exhibit explicitly lexical representations, they may take complementary
146 roles consistent with the dual stream pathway model (Gow, 2012; Hickok and Poeppel, 2007).
147 Thus, the type of task (e.g., whether subjects are asked to make a phonological or semantic
148 judgment) may influence the relative dominance between SMG and MTG activities during
149 speech recognition.

150 Supporting the idea of broad cortical regions contributing to individual differences in
151 speech processing, fMRI studies showed that SNR changes alter the level of neural activities
152 across frontal, central, and temporo-parietal regions (Du et al., 2016; Vaden et al., 2015; Wong
153 et al., 2009; Zekveld et al., 2006), while Du et al. (2016) reported the correlation between
154 activities in fronto-central regions and speech recognition performance. However, these
155 correlations could reflect earlier variation in speech unmasking – if auditory/attentional
156 unmasking mechanisms are less efficient, then they could lead to differences in how strongly
157 later regions (IFG, SMG, etc.) must work to recognize words or complete the task. Thus, it is
158 crucial to evaluate both mechanisms simultaneously to isolate a potential role for later
159 processes.

160 In addition to testing the loci of activities, it is also important to test the relative timing of
161 activity in these pathways during speech processing. Functionally, studies using eye-
162 movements in the Visual World Paradigm (VWP) have extensively characterized the time
163 course of word recognition in both quiet (Allopenna et al., 1998; Dahan and Gareth Gaskell,
164 2007; Magnuson et al., 2007) and under challenging conditions such as noise or signal
165 degradation (Ben-David et al., 2011; Brouwer and Bradlow, 2016; Huettig and Altmann, 2005;
166 McMurray et al., 2017; McQueen and Huettig, 2012). Most VWP data show that, in quiet, the
167 maximum lexical competition occurs at ~400 ms after the onset. These studies also report
168 delayed processing under challenging conditions, but such delays do not exceed 250 ms even
169 under the most severe degradation. This timing information can guide us when we interpret the
170 functional implication of neural activity within certain regions; if the latency of neural activity is
171 larger than ~400 ms, such activity may be less likely related to the online word recognition.

172 However, the timing of cortical activity within each pathway and the way this may be
173 moderated by challenging listening conditions are largely unknown. This is true both within
174 isolated regions (e.g., SMG or MTG) and across activation of broader regions (e.g., frontal
175 lobe). This is because most of the work on speech in noise perception has been conducted with

176 fMRI (Du et al., 2014, 2016; Wong et al., 2009; Wong et al., 2008) which has a poor temporal
177 resolution. One study has examined evoked responses across speech processing regions using
178 source localized EEG (Bidelman and Howell, 2016). This suggests an early response at roughly
179 100 ms post-stimulus in IFG. However, this study used non-sense syllables that cannot reveal
180 lexical processes.

181 ***The present study.*** The central aim of this study is to investigate the simultaneous
182 contributions of both speech unmasking and speech recognition processes to the individual
183 differences in SiN understanding. Our main question is whether the early-stage speech
184 *unmasking* and later-stage *recognition* processes independently predict SiN performance, or
185 whether the latter variable is dependent on the former. We attempted to answer this question
186 through the combination of within- and across-subject analyses using both sensor- and source-
187 space evoked responses in an EEG paradigm.

188 Subjects performed a SiN noise task in which they heard isolated consonant-vowel-
189 consonant (CVC) English words and selected which of four orthographically presented words
190 matched the auditory stimuli. Noise began 1 second before the speech. Two noise conditions
191 were used: a low SNR (-3 dB) or hard condition, and high SNR (+3 dB) or easy condition. EEG
192 was recorded from 64 electrodes while subjects performed this task, using both source- and
193 sensor-space analyses to quantify cortical activity.

194 Speech unmasking and speech recognition can be distinguished by the 1) timing and 2)
195 regional differences of neural activities evoked by both noise and target speech. Thus, we used
196 a trial structure comprised of clearly separated events (i.e., fixed onsets of background noise
197 and target speech) while observing time-locked neural responses to such events with
198 electroencephalography (EEG). The degree of speech unmasking was quantified as the
199 amplitude ratio of evoked responses to the onsets of noise and target speech measured at a
200 vertex scalp electrode (the key scalp location for evoked responses from early auditory
201 processes), henceforth referred to as *internal SNR*. Although the concept of *internal SNR* is

202 similar to the attentional modulation index (Dai and Shinn-Cunningham, 2016; O'Sullivan et al.,
203 2019), we named the index rather phenomenologically to avoid limiting the mechanism
204 underlying the index to selective attention.

205 The effectiveness of later speech *recognition* processing was quantified by measuring
206 the amplitude of evoked responses within target cortical regions. As we described, prior work
207 has implicated a number of such regions. However, it is unclear which may be relevant for our
208 specific task. Thus, we take a data-driven approach by first asking which post-auditory regions
209 show *greater* activity in the *higher* SNR (easier) listening condition. This would be suggestive of
210 a region that conducts downstream analyses once the target speech is unmasked. We looked
211 into two regions-of-interest (ROIs): left SMG and IFG. As reviewed above, those regions
212 activate more strongly for speech than pseudo-speech sounds; we did not consider MTG (the
213 ventral lexicon) as our task was single CVC word identification (matched to an orthographic
214 response). In this task, phonological discrimination (indicating the dorsal stream), rather than
215 semantic processing, is essential.

216 Having quantified the contributions of each pathway, we then conducted both timing and
217 individual differences analyses to determine their relative contribution to SiN performance.

218

219 **2. Material and Methods**

220 **2.1 Participants**

221 Twenty-six subjects between 19 and 31 years of age (mean = 22.42 years, SD = 2.97
222 years; median = 21.5 years; 8 (31%) male) were recruited from a population of students at the
223 University of Iowa. All subjects were native speakers of American English, with normal hearing
224 thresholds no worse than 20 dB HL at any frequency, tested in octaves from 250 to 8000 Hz.
225 Written informed consent was obtained, and all work has been carried out in accordance with
226 the Code of Ethics of the World Medical Association (Declaration of Helsinki). All study
227 procedures were reviewed and approved by the University of Iowa Institutional Review Board.

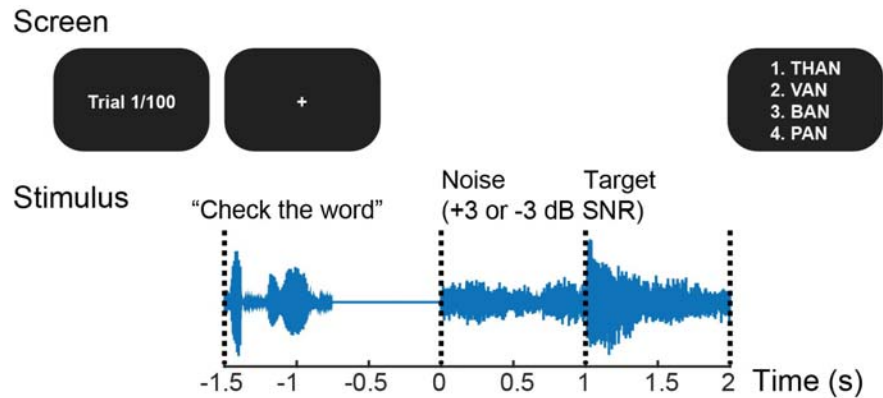
228

229 **2.2 Task design and procedures**

230 We aimed to simultaneously measure SiN performance and cortical neural activity in a
231 short (15 minute) experimental session. Sessions were kept short to avoid confounding
232 individual differences in irrelevant psychological factors – fatigue, level or engagement –with
233 individual differences in performance and processing.

234 Each trial (**Figure 1**) began with the presentation of a fixation cross ('+') on the screen.
235 Listeners were asked to fix their gaze on this throughout the trial to minimize eye-movement
236 artifacts. Next, they heard the cue phrase “check the word.” This enabled listeners to predict the
237 timing of next acoustic event (the noise onset). After a 700 ms of silence, the multi-talker babble
238 noise began and continued for 2 seconds. One second after the noise onset, the target word
239 was heard. Finally, 100 ms after the composite auditory stimulus (noise + word) offset, four
240 written choices appeared on the screen. The response options differed either in the initial or the
241 final consonant (e.g., for target word *ban*, options were *than*, *van*, *ban*, and *pan*; for target word
242 *hiss*, options included *hit*, *hip*, *hiss*, *hitch*). Subjects pressed a button on a keypad to indicate
243 their choice and no feedback was given. The next trial began 1 second after the button press.

244 This trial structure was intended to minimize visual, pre-motor, and motor artifacts during
245 the time of interest surrounding the auditory stimuli. The timing and intervals of auditory stimuli
246 (i.e., cue phrase, noise, and target) were intended to derive well-distinct cortical evoked
247 responses to the onsets of background noise and target word.



249 **Figure 1.** Trial and stimulus structure. Every trial starts with the cue phrase “check the word.” A target
250 word starts 1 second after the noise onset. Four choices are given after the word ends; subjects select
251 the correct answer with a keypad. No feedback is given. The noise level is manipulated to create high (+3
252 dB) and low (-3 dB) SNR conditions. Subjects complete 50 trials for each condition.
253

254 Since we were particularly interested in SMG and IFG regions that are involved in
255 phonological and lexical processing (Gow, 2012; Hickok and Poeppel, 2007), we used naturally
256 spoken words, rather than non-sense speech tokens used by prior EEG studies (Bidelman and
257 Howell, 2016; Parbery-Clark et al., 2009). Target words consisted of hundred monosyllabic CVC
258 words from the California Consonant Test (CCT) (Owens and Schubert, 1977), spoken by a
259 male speaker with a General American accent.

260 Target words were always presented at 65 dB SPL. In each trial, the RMS level of noise
261 was chosen randomly between 68 and 62 dB SPL to yield either -3 or +3dB SNR (referred to as
262 “low SNR” and “high SNR,” respectively). Fifty words were presented at each SNR. -3 dB SNR
263 was chosen from pilot experiments to emulate a condition yielding a mid-point performance
264 (~65% correct) in the possible accuracy range (i.e., 25 – 100%), at which listening effort and
265 individual differences in performance may be maximized (Ohlenforst et al., 2017). Thus, the SiN
266 performance at -3 dB SNR condition was used for the later correlational analysis in this study.
267 +3 dB SNR was chosen to emulate a less noisy condition from which the downstream speech
268 recognition process will be measured for the correlational analysis.

269 The task was implemented using the Psychtoolbox 3 package (Brainard, 1997; Pelli,
270 1997) for Matlab (R2016b, The Mathworks). Participants were tested a sound-treated,
271 electrically shielded booth with a single loudspeaker (model #LOFT40, JBL) positioned at a 0°
272 azimuth angle at a distance of 1.2 m. A computer monitor was located 0.5m in front of the
273 subject at eye level. The auditory stimuli were presented at the same levels for all subjects.

274

275 **2.3 EEG acquisition and preprocessing**

276 Scalp electrical activity (EEG) was recorded during the SiN task using the BioSemi
277 ActiveTwo system at a 2048 Hz sampling rate. Sixty-four active electrodes were placed
278 according to the international 10-20 configuration. Trigger signals were sent from Matlab
279 (R2016b, The Mathworks) to the ActiView acquisition software (BioSemi). The recorded EEG
280 data from each channel were bandpass filtered from 1 to 50 Hz using a 2048-point FIR filter.
281 Epochs were extracted from -500 ms to 3 s relative to stimulus onset. After baseline correction
282 using the average voltage between -200 and 0 ms, epochs were down-sampled to 256 Hz.

283 Since we were interested in the speech-evoked responses from frontal brain regions, we
284 opted for a non-modifying approach to eye blink rejection: Trials that were contaminated by an
285 eye blink artifact were rejected based on the voltage value of the Fp1 electrode (bandpass
286 filtered between 1 and 20 Hz). Rejection thresholds for eye blink artifacts were chosen
287 individually for each subject, and separately for the noise and target word periods. After
288 rejecting bad trials, averages for each electrode were calculated for the two conditions to extract
289 evoked potentials. For analysis of speech-evoked responses, we repeated baseline correction
290 using the average signal in the 300 ms preceding the word onset.

291

292 **2.4 Sensor-space analysis**

293 We performed traditional sensor-space ERP analysis to investigate the effect of acoustic
294 SNR on AC representation of noise and speech, and its individual differences. The other

295 purpose of sensor-space analysis was to ensure the quality of data in the more familiar form
296 before running source-space analyses.

297 Cortical evoked responses time-locked to the target and noise were examined and
298 compared between high- and low-SNR conditions. EEG data were bandpass filtered from 2 to 7
299 Hz to capture auditory N1 and P2 components that fall into 3 – 5 Hz bands. Our first, analysis
300 examined the N1 and P2 peak amplitude obtained from the frontal-central channels (C1, C2,
301 FC1, FC2, FCz, Cz) and compared these measures between two SNR conditions. Both auditory
302 N1 and P2 components were obtained at around 200 ms after the noise onset and at about 250
303 ms after the target word onset. After comparing the amplitude of each auditory component, we
304 also examined the ERP envelopes by applying the Hilbert transform to the bandpass-filtered
305 ERPs and taking the absolute value to effectively represent overall magnitude of both N1 and
306 P2 ERP components. Then, “internal SNR” was defined as the ratio of target word-evoked ERP
307 envelope peaks to noise-evoked ERP envelope peaks magnitude in dB scale (Equation 1). We
308 computed this index expecting to quantify a “neural” form of an individual’s speech unmasking
309 ability. The internal SNR is different for each subject, and is separate from the fixed external, or
310 acoustic, SNR (here, ± 3 dB).

311

312
$$(1) \text{ Internal SNR} = 20 \log_{10} \frac{\text{Target-evoked potential}}{\text{Noise-evoked potential}}$$

313

314 Mean activity levels were jackknifed prior to testing to assess the variance with clean
315 ERP waveforms. In this approach, the relevant neural factors were computed for all subjects but
316 one. This was repeated leaving out each subject in turn. The resulting statistics were adjusted
317 for jackknifing to reflect the fact that each data point reflects N-1 subjects (Luck, 2014).

318

319 **2.5 Source analysis**

320 The source-space analysis was based on minimum norm estimation (Gramfort et al.,
321 2013; Gramfort et al., 2014) as a form of multiple sparse priors (Friston et al., 2008). After co-
322 registration of average electrode positions to the reconstructed average head model MRI, the
323 forward solution (a linear operator that transforms source-space signals to sensor space) was
324 computed using a single-compartment boundary-element model (Hämäläinen, 1989). The
325 cortical current distribution was estimated assuming that the orientation of the source is
326 perpendicular to the cortical mesh. Cross-channel EEG-noise covariance, computed for each
327 subject, was used to calculate the inverse operators. A noise-normalization procedure was used
328 to obtain dynamic statistical parametric maps (dSPMs) as z-scores (Dale et al., 2000). The
329 inverse solution estimated the source-space time courses of event-related activity at each of
330 10,242 cortical voxels per hemisphere.

331 Following the whole-brain source estimation, we extracted representative source time
332 courses from regions of interests (ROIs). In the present study, two predetermined ROIs were
333 used: (1) left SMG, and (2) left pars opercularis and pars triangularis of IFG. Destrieux Atlas of
334 cortical parcellation (Fischl et al., 2004) was used to predetermine ROIs anatomically.

335 Since we did not have individual structural MRI head models, it was not ideal to take the
336 summed activity (mean or median) for all the voxels within ROIs. This is because individual
337 difference in functional and anatomical structure of the brain may result in spatial blurring since
338 current densities across adjacent voxels can overlap each other. Instead, representative voxels
339 were identified for each ROI, for each SNR condition. We used a combination of previously-
340 described methods to select voxels of interest that were used in fMRI studies (Tong et al.,
341 2016). The voxel selection was performed by a two-step process. First, we selected voxels that
342 exhibit greater-than-median amplitude in either (high **or** low SNR condition) condition. Second,
343 cross-correlation coefficients for ERP time courses across all remaining voxels in an ROI were
344 calculated across time, and then the mean coefficient was calculated for each voxel. The most
345 representative voxel was defined as having the maximum mean correlation coefficient, while

346 also being above threshold at two or more continuous timepoints based on voxel's p -value, as
347 determined using one-sample t -tests (Tong et al., 2016).

348 For the downstream statistical analyses, temporal envelopes were extracted from the
349 within-ROI source time courses. The source time course envelopes at different ROIs were
350 based on slightly different frequency ranges: 2-7 Hz for SMG and 1-5 Hz for IFG. This was to
351 reflect the differences between temporal and frontal lobes in their dominant neural oscillations
352 which create evoked activities through phase coherence (Giraud and Poeppel, 2012).

353

354 **2.6 Statistical approaches**

355 Once the temporal envelope of the source time course from the most representative
356 voxel was obtained for each SNR condition, mean activity levels were compared between the
357 two SNR conditions using paired t -tests. Here, we also used a jackknifing approach, and test
358 statistics were adjusted to account for the fact that each data-point represents $N-1$ participants.
359 Finally, to identify timepoints that showed a significant difference between SNR conditions while
360 addressing multiple comparison problem, the cluster-based permutation tests were conducted
361 (Maris and Oostenveld, 2007).

362 In order to identify predictors of SiN performance, sensor and source space indices of
363 activity were used in correlation/regression analysis with SiN performance (accuracy) as the
364 dependent variable, and the peak magnitudes of the ERP envelopes from ROIs, and internal
365 SNR as the predictor variables. The peak magnitudes of the ERP envelopes were obtained over
366 timepoints that showed a significant difference between high and low SNR conditions identified
367 in the ROI-based source analysis described above. After calculating the correlation between
368 SiN performance and these predictors, a joint contribution was tested using linear regression
369 analysis to simultaneously examine bottom-up and compensatory related SiN performance to
370 three factors.

371

372 **3. Results**

373 Our analysis started by examining the effect of SNR on task performance (accuracy and
374 reaction time). This was intended to document that noise manipulation had the expected effect.
375 Next, we examined the effect of SNR on both the magnitude and timing of neural activity. This
376 was done first in the sensor-space, using the auditory N1/P2 components and other measures
377 to examine primary auditory pathways. Next, this was done in the source-space to examine the
378 compensatory role of IFG, to evaluate whether IFG effects were early enough to play a role in
379 speech perception, and to test hypotheses about SMG. Finally, we turn to our primary analysis:
380 a regression testing the unique contributions of each pathway to individual differences in SiN
381 performance. Original raw and processed data of the present study are available at Mendeley
382 Data (<http://dx.doi.org/10.17632/jvythkz5y.1>).

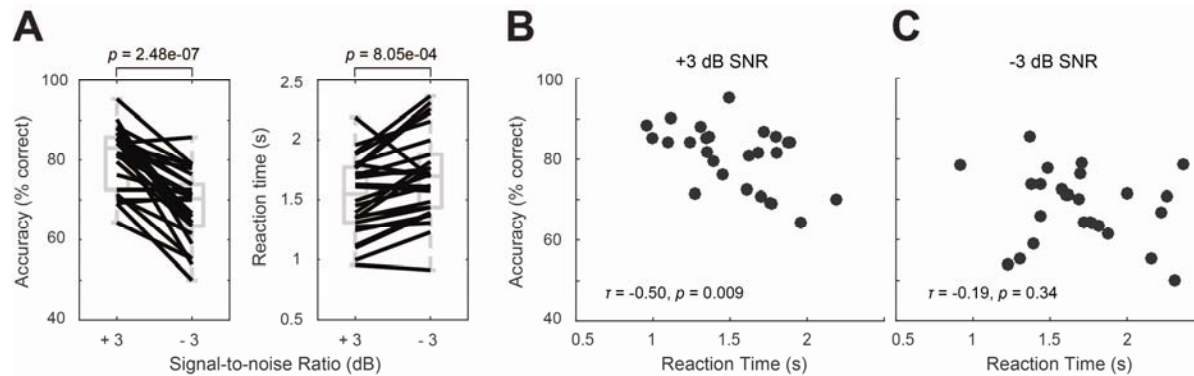
383

384 **3.1 SiN performance**

385 There was a large variance in performance among participants. This was observed in
386 both the high SNR condition (accuracy: mean = 80.64%, SD = 7.81%, median = 83.01%;
387 reaction time: mean = 1.53 s, SD = 0.32 s, median = 1.55 s) and the low SNR condition
388 (accuracy: mean = 68.21%, SD = 8.92%, median = 70.37%; reaction time: mean = 1.70 s, SD =
389 0.36 s, median = 1.69 s). There was a significant effect of SNR on both accuracy ($t(25) = 6.99$, p
390 < 0.001 , paired t-test) and reaction time ($t(25) = -3.81$, $p < 0.001$, paired t-test) (**Figure 2A**).
391 Reaction time and accuracy were correlated in the high SNR condition (**Figure 2B**, $r = -0.50$, p
392 $= 0.009$), but not in the low SNR condition (**Figure 2C**, $r = -0.19$, $p = 0.34$).

393 As a whole, these results validate that the SNR manipulation was sufficient to create
394 differences in speech perception. The negative correlation between the accuracy and reaction

395 time may demonstrate the redundancy between individual differences of accuracy and difficulty.



396

397 **Figure 2.** Behavioral results. **A.** Summary of behavioral performance for the two conditions (+3 and -3 dB
398 SNR). Boxes denote the 25th – 75th percentile range; the horizontal bars in the center denote the median;
399 the ranges are indicated by vertical dashed lines. Solid lines connect points for the same subject in
400 different conditions. **B.** Average accuracy as a function of reaction time in +3 dB SNR condition. **C.**
401 Average accuracy and reaction time in -3 dB SNR condition.
402

403

3.2 Auditory N1-P2 components in the sensor space

404

We next examined the N1 and P2 peak amplitude for both noise- and target word-
405 evoked responses. This was done to estimate the contribution of primary auditory pathways and
406 the effect of noise. These measured were compared between high- and low-SNR conditions.

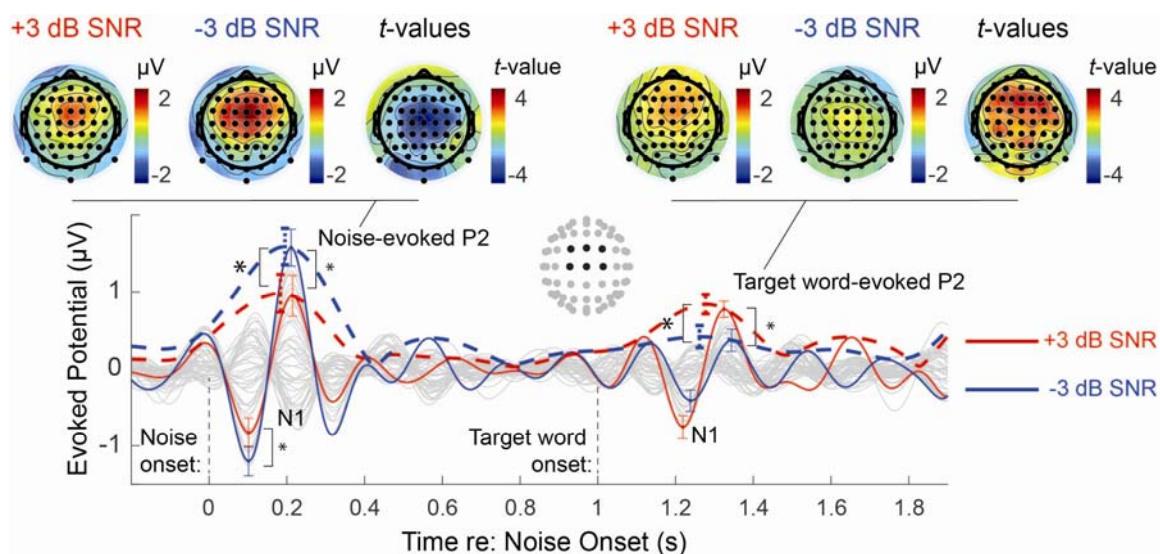
407

We started with the noise-evoked responses (i.e., peak ERP components arising ~0.2
408 seconds following the noise onset at 0 s). Here, both N1 ($t(25) = -2.64$, $p = 0.014$) and P2
409 amplitudes ($t(25) = -3.82$, $p < 0.001$) were greater in the low SNR condition. In contrast, target-
410 word evoked responses (i.e., peak ERP components arising ~0.25 seconds following the target-
411 word onset at 1 s) showed the opposite pattern with stronger N1 ($t(25) = 1.87$, $p = 0.073$) and
412 P2 amplitude ($t(25) = 2.41$, $p = 0.023$) in the high-SNR condition (**Figure 3**).

413

We also examined the ERP envelopes in the time region of both the auditory N1 and P2.
414 Again in the noise-evoked time region, we saw a greater response in the low SNR condition
415 ($t(25) = -3.95$, $p < 0.001$). However, also as in the N1/P2 analysis we saw larger word-evoked
416 response in the high SNR condition ($t(25) = 2.37$, $p = 0.026$) (**Figure 3**). Then we calculated the
417 internal SNR, the magnitude ratio of the noise and target-word related ERP envelopes, and saw
418 a greater internal SNR in the high SNR condition ($t(25) = 2.53$, $p = 0.018$).

419 The topographical layout of these effects is shown in **Figure 3** (top panels) at the time of
420 the peak P2 (which showed a significant overall effect of noise level) for both noise- and word-
421 evoked response. T-values from paired t-tests on peak P2 amplitudes at all electrodes between
422 SNR conditions were also represented in topographies. These show a broad-based effect that is
423 roughly centered at frontal-central channels for both the noise onset and speech onset (though
424 patterns for the speech evoked P2 were more distributed along the scalp). This justifies our use
425 of frontal-central channels for sensor-space analyses. The significant difference in auditory ERP
426 envelopes and in the N1/P2 components according to the noise level supports our use of the
427 internal SNR as an index of individual ability to modulate representations of target speech
428 relative to noise in the regression analysis.



429

430 **Figure 3.** Sensor-level event-related potential (ERP) with the topographical layout and cortical maps. The
431 time course of the auditory ERP and its envelope, with the standard error of the mean (± 1 SEM) at the
432 peak amplitude (red color: +3 dB SNR, blue color: -3 dB SNR). An asterisk shows a significant difference
433 in the amplitude between +3 and -3 dB SNR conditions (paired t -test). Top panels show peak P2
434 amplitudes of all electrodes in topographical layouts. The t -values from paired t -tests between two SNR
435 conditions are also shown as topographies.

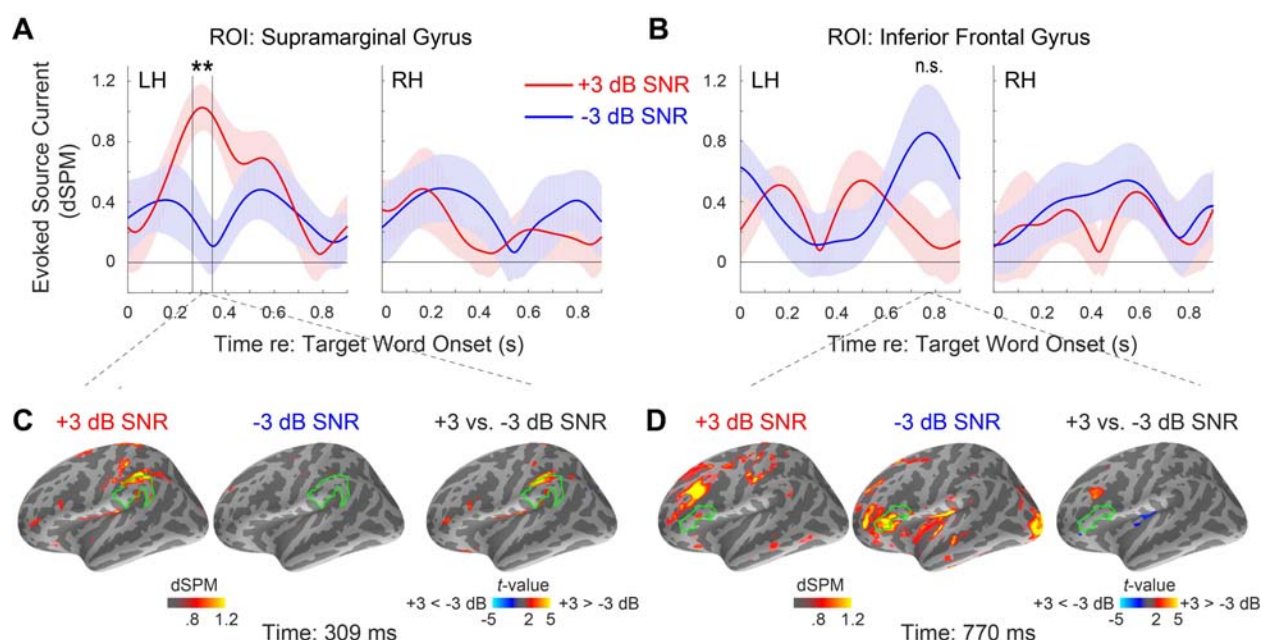
436

437 **3.3 The effect of SNR on cortical activity**

438 Next, we conducted parallel analyses in source space to assess cortical activity through
439 speech processing regions. We converted sensor-space EEG signals to whole-brain source

440 time courses to localize the effects of SNR on evoked responses within targeted ROIs. Within
 441 left SMG, the cluster-based permutation test (Maris and Oostenveld, 2007) revealed that the
 442 high SNR condition evokes significantly greater activity than the low SNR condition from 270 to
 443 340 ms ($p = 0.0020$) (**Figure 4A left**). High-SNR peak amplitude is found at 309 ms. Such a
 444 significant SNR effect was not found in the left IFG. The maximum magnitude of the grand
 445 average low-SNR evoked response (dSPM) is at 770 ms (**Figure 4B left**). Source time courses
 446 in the right SMG and IFG are shown in **Figure 4A and B** (the right panels) for visual
 447 comparisons.

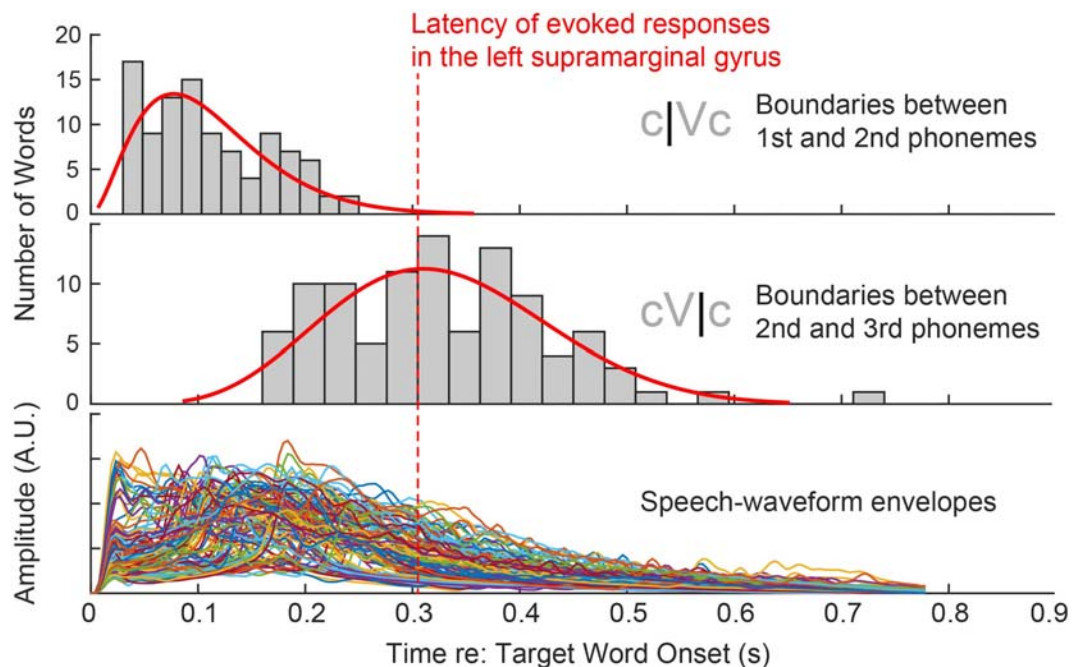
448 Single time-point evoked-current estimates are shown for the peak SMG activity time
 449 (i.e., 309ms, **Figure 4C**) and IFG peak time (770ms, **Figure 4D**) on the whole left-hemisphere
 450 cortical surface. At 309ms, post-hoc paired t -tests on all the left-hemisphere voxels reveal an
 451 area near left SMG that shows greater evoked responses in the high SNR condition. This
 452 supports a significant role for SMG in SiN processing in this task. In contrast, at 770ms, no
 453 voxel was found within IFG that shows significant differences between high- vs. low-SNR
 454 conditions. This confirms our timecourse analyses of IFG, suggesting it does not play a large
 455 role and that the trend that was observed is not broadly seen across voxels.



456

457 **Figure 4.** Region-of-interest (ROI) based source analysis. **A and B.** The time course of the event-related
458 potential (ERP) envelope, with the standard error of the mean (± 1 SEM), obtained at representative
459 voxels for two ROIs in the left hemisphere ("LH"): supramarginal gyrus (SMG), and the pars opercularis
460 and triangularis of the inferior frontal gyrus (IFG), respectively, in each SNR condition (red color: +3 dB
461 SNR, blue color: -3 dB SNR). The ERP envelope time courses in the corresponding regions in the right
462 hemisphere ("RH") are also shown for visual comparisons. Asterisks show the timing of a significant
463 difference between +3 and -3 dB SNR conditions (cluster-based permutation test, $p = 0.0020$) at the left
464 SMG. **C and D.** Whole-brain maps showing statistical contrasts (t -values obtained from post-hoc paired t -
465 tests between the two SNR conditions) of source activation at each voxel at the peak timepoint of the
466 grand-average source time course of each ROI.
467

468 To demonstrate the timing of SMG activity compared to the timing of phonological
469 events, the webMAUS (Kisler et al., 2017) was used to identify the boundaries between the first
470 and second, and the second and third phonemes in each of the 100 stimuli. A histogram of
471 these acoustic time points is shown in **Figure 5**. This confirms that the peak of evoked
472 activation in SMG (i.e., ~309 ms, denoted by a dashed vertical line) occurs within the time
473 course of target words before the final phoneme is presented.



474 **Figure 5.** Top and second panel show histograms of boundaries between phonemes of each stimulus.
475 The third panel shows superimposed temporal envelopes extracted from waveforms of the 100 words.
476
477

478 **3.4 Individual differences in internal SNR predict SiN performance**

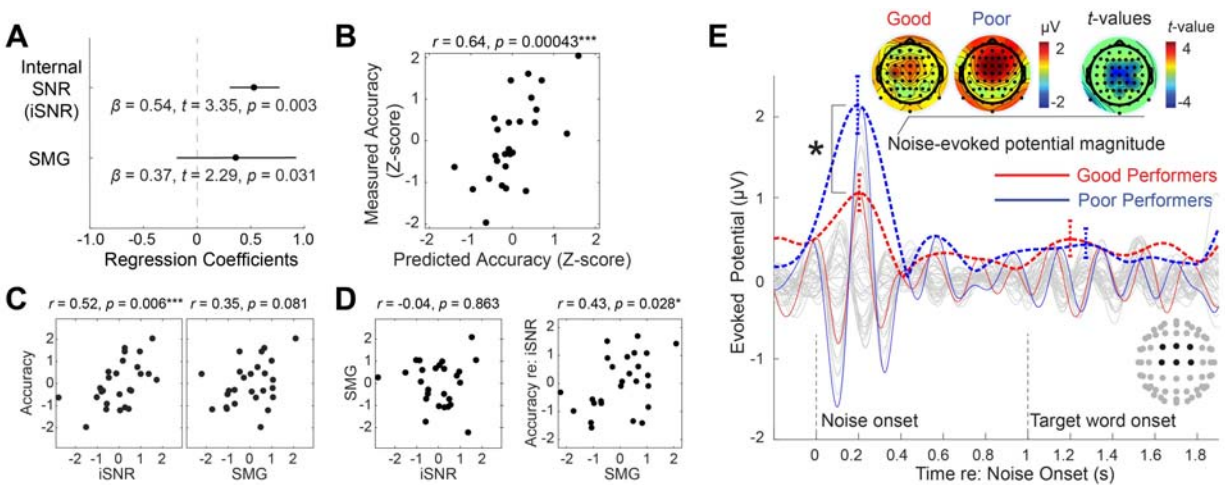
479 To address our primary research question, which was to evaluate the simultaneous
480 contribution of speech unmasking and recognition processes to SiN performance, we conducted
481 a linear regression analysis in which internal SNR and SMG activation were used as
482 independent variables. SiN performance in low SNR condition, which showed larger variance
483 (high SNR SD = 7.81%; low SNR: SD = 8.92%), was used as the dependent variable. We
484 extracted the internal SNR from the low SNR and the SMG activity from the high SNR condition,
485 as we expected that the internal SNR captures how well listeners unmask speech from the
486 noisy background while the SMG activity reflects the processing of relatively clean speech
487 signal. As expected, those two metrics extracted from different trials did not show a correlation
488 ($r = -0.04$, $p = 0.863$, the left panel of **Figure 6D**). Both internal SNR ($t(23) = 3.35$, $p = 0.003$)
489 and SMG activity ($t(23) = 2.29$, $p = 0.031$) were significant predictors of SiN performance
490 (**Figure 6A**). The linear combination of those predictors accounted for a large proportion of the
491 variance ($r = 0.64$, $p = 0.00043$, **Figure 6B**).

492 **Figure 6C and D** show results from post-hoc correlational analyses. Internal SNR
493 showed a significant correlation with accuracy, while SMG activation did not (despite its
494 significant contribution to the model). There was no correlation between internal SNR and SMG
495 activation, as described above. A semi-partial correlation between SMG activation and the
496 residual of accuracy after regressing out internal SNR was significant, which confirmed that the
497 SMG activation accounted for an extra amount of variance in SiN performance (the right panel
498 of **Figure 6D**). This suggests that in order to identify the contribution of downstream recognition
499 areas like SMG, models must account for the contribution of earlier upstream speech-
500 unmasking processes.

501 To visualize the contribution of internal SNR to SiN performance, **Figure 6E** showed
502 evoked response differences between good and poor performers (based on a median split on
503 the low-SNR condition accuracy). This reveals dramatic differences in the magnitude of noise
504 onset-related potentials: despite the same physical noise level for each group, good performers

505 exhibited less strong evoked response to noise onset, measured by the envelope peak
 506 magnitude within N1-P2 time range in the frontal-central channels (two-sample t -test $t(24) = -$
 507 2.60, $p = 0.016$). In contrast, the word-evoked ERP envelope did not show a significant
 508 difference between the two groups ($t(24) = 0.21$, $p = 0.84$, two-sample t -test). This suggests that
 509 the neural mechanism underlying the internal SNR variance is the suppression of noise (rather
 510 than the enhancement of the target).

511



512

513 **Figure 6.** Individual differences in speech-in-noise processing. **A.** Regression coefficients and their
 514 standard errors. **B.** A scatter plot showing the relationship between predicted and measured accuracy in -
 515 3 dB SNR condition. **C.** Post-hoc correlation analyses: Raw correlations between each independent
 516 variable and the dependent variable. **D.** Left: Relationship between independent variables shows no
 517 correlation between internal SNR and evoked source current at the left supramarginal gyrus (SMG).
 518 Right: Semi-partial correlation between SMG evoked source current and the residual of accuracy after
 519 regressing out internal SNR. **E.** The time course of the auditory event-related potential and its envelope,
 520 with the standard error of the mean (± 1 SEM) at the peak magnitude in -3 dB SNR condition (red color:
 521 good performers, blue color: poor performers). An asterisk shows a significant difference in the magnitude
 522 between two groups (two-sample t -test).

523

524 4. Discussion

525 We investigated neural correlates of SiN performance in young normal hearing adults.
 526 Previous correlational studies focused on the contributions of either acoustic encoding fidelity
 527 (Anderson and Kraus, 2010; Anderson et al., 2013; Holmes and Griffiths, 2019; Hornickel et al.,
 528 2009; Liberman et al., 2016; Parbery-Clark et al., 2009; Song et al., 2011) or the degree of

529 speech/language network recruitment (Du et al., 2016) to the SiN performance. However, the
530 relative importance of each process has remained unclear. We showed that 1) how well the
531 listener suppresses background noise *before* hearing the target speech and 2) how strongly the
532 listener recruits temporo-parietal network *while* the speech signal is received contribute to the
533 SiN performance independently. Combining those two factors explained about 40% of the
534 variance in SiN performance.

535 Our results have both theoretical and clinical implications. Theoretically, our individual
536 difference approach revealed at least two neural subsystems involve during SiN processing:
537 sensory gain control and post-auditory speech recognition processing. Clinically, our results
538 suggest that a relatively short (~15 minutes) SiN-EEG paradigm can assess crucial neural
539 processes for SiN understanding.

540 ***Internal SNR: A measure of pre-speech processing for speech unmasking.*** The first
541 among the two crucial processes – how well the listener suppresses background noise – was
542 indexed as “internal SNR,” the ratio of noise- to target word-evoked cortical responses. This
543 process can be understood as pre-target cortical activity, appearing as an enhanced neural
544 representation of the target sound (the speech) and suppressed neural representation of
545 ignored stimuli (the noise).

546 ***What is the source of variation in the internal SNR?*** Such responses could reflect
547 auditory selective attention, which shows a similar pattern in previous studies (Hillyard et al.,
548 1973; Hillyard et al., 1998; Mesgarani and Chang, 2012). In the present study, good performers
549 showed significantly weaker noise-evoked responses at frontal-central channels (around Cz),
550 compared with poor performers, approximately 200 ms after the noise onset (**Figure 6E**).
551 Decreased auditory responses to background noise in good performers are compatible with the
552 presence of a sensory gain control mechanism (Hillyard et al., 1998) which may happen in
553 multiple sub-regions in STP and posterior STG. The variation in the sensory gain control may
554 originate from multiple factors. It may reflect the acuity of encoding spectro-temporal acoustic

555 cues from speech and noise, or grouping of such acoustic cues for auditory object formation
556 (Moore, 1990; Shamma et al., 2013; Teki et al., 2011). How robustly the low-frequency neural
557 oscillations (e.g., theta and delta) are phase-locked to the acoustic temporal structure of the
558 stimuli (Etard and Reichenbach, 2019) may also contribute to the variation, as the neural phase-
559 locking relies on the encoding of acoustic cues (Ding et al., 2014) and the prediction of temporal
560 structure in speech rhythm (Ding et al., 2016). Since our experiment provided fixed timing of
561 noise and target word onsets, the neural phase-locking based on predicted timing (Arnal and
562 Giraud, 2012) could occur and contribute to the internal SNR. The variation may also reflect
563 endogenous mechanisms for active suppression of background sounds along with neural
564 enhancement of foreground sounds (Shinn-Cunningham and Best, 2008). It was not our goal to
565 disentangle the sources of variation in sensory gain control. Rather, we aimed to quantify the
566 effectiveness of sensory gain control by our unique trial structure that enables clear distinction
567 of evoked responses to noise and target speech, and to test how the internal SNR predicts later
568 speech processes and behavioral accuracy. In this regard, we found a significant correlation
569 between accuracy and the relative magnitude of word- and noise-evoked potentials.

570 ***Evoked amplitude in SMG: the neural marker of effective and prompt lexical***
571 ***processing.*** While the computation of internal SNR was pre-specified, we had an open plan for
572 extracting a representative neural factor to capture post-auditory speech recognition. To explore
573 such neural markers, we added a 6-dB higher SNR condition and asked which region or regions
574 showed increased activity within a reasonable (200 – 500ms) time range. We investigated two
575 ROIs: left SMG and left IFG. As left SMG showed increased evoked response to target speech
576 in the less noisy condition at ~300 ms after the target onset, the peak evoked amplitude in left
577 SMG measured in the high SNR condition was used as the second independent variable in the
578 regression analysis.

579 ***Functional interpretation of SMG activity.*** Previous studies have suggested that
580 spoken-word recognition occurs via a process of dynamic lexical competition as speech unfolds

581 over time. The VWP studies reported that, for many words, this competition maximizes around
582 3-400 ms after word onset (Farris-Trimble and McMurray, 2013; Huettig and Altmann, 2005). In
583 significantly challenging conditions (high noise), however, lexical processing can be delayed
584 about 250 ms until most of the word has been heard (Farris-Trimble et al., 2014; McMurray et
585 al., 2017), which may minimize competition. The latency of SMG activity that lied between the
586 second and the third phonemes (see **Figure 5**) in the high SNR condition aligns well with the
587 timing of lexical competition found from the VWP studies, which may suggest that the SMG
588 activity makes a neural substrate of immediate lexical access (Farris-Trimble et al., 2014;
589 McMurray et al., 2017), consistent with Gow (2012). This immediacy was observed when
590 speech sounds were relatively clean (high SNR), and it does not appear in previous EEG
591 studies using non-word synthesized phonemes (Bidelman and Dexter, 2015; Bidelman and
592 Howell, 2016).

593 After the contribution of speech unmasking (i.e., internal SNR) is regressed out, the
594 SMG evoked amplitude in the cleaner condition predicted the residual of SiN performance (the
595 right panel of **Figure 6D**). This indicates that changes in SMG activity may be an independent
596 factor predicting speech recognition performance, rather than the outcome of pre-speech
597 sensory gain control processing.

598 ***Limitation of the current study.*** In the present paper, we exhibited evoked responses
599 only. Although our results demonstrated how this simple and traditional EEG analysis
600 successfully predicted SiN performance, future studies may pursue further understanding of SiN
601 mechanisms by adopting extended analyses such as induced oscillation (e.g., Choi et al.
602 (2020)) and connectivity analyses.

603 Our correlational result is limited to young normal hearing listeners where variance does
604 not come from hearing deterioration and aging. This study does not predict how much of the
605 variance can be explained by the combination of internal SNR and SMG activity in the

606 population with larger age ranges. For example, Tune et al. (2020) reported that no correlation
607 is found between neural factors and behavioral success in a large cohort of aged listeners.

608 We chose -3 dB as the main SNR from which the dependent variable for the regression
609 analysis was extracted. Although we claim that the -3 dB SNR provides the most representative
610 condition where the individual difference in performance is maximized, we do not claim that our
611 main findings can be generalized to other SNR conditions.

612 ***Methodological advances and justifications for source time course analysis.*** Our
613 approach to identifying a single voxel within an ROI deserves a particular discussion.
614 Identification of the representative voxel of an ROI is a problem common to EEG source
615 analysis, fMRI, and other functional brain imaging studies. Many relevant neuroimaging analysis
616 approaches have been described, including univariate, multivariate, and machine learning;
617 however, most of these are intended for the identification of regions of interest or functional
618 connections from a whole-brain map. Drawbacks of this type of whole-brain analysis include the
619 need for strict multiple comparisons correction and, therefore, decreased statistical power.
620 Using strong a priori hypotheses to generate regions of interest allowed us to circumvent these
621 issues, but still requires identification of representative voxels within our regions of interest.
622 Favored approaches generally require identification of peak activity within an ROI (Tong et al.,
623 2016). However, to avoid the assumption that choosing peak activity implies, we opted instead
624 to choose the voxel that has the maximum average correlation to every other voxel within the
625 ROI. In the present study, we chose not to constrain the location of the voxel of interest within
626 an ROI for each condition. Because our anatomic resolution is unlikely to be at the voxel level,
627 we elected to choose a different representative voxel for each condition, unconstrained by the
628 location of the representative voxel from other conditions.

629

630 **5. Conclusion**

631 We found that better speech unmasking in good performers modulated the ratio of
632 cortical evoked responses to the background noise and target sound, which effectively changed
633 SNR internally, resulting in better performance. We also found that clean, intelligible speech
634 elicits early processing at SMG, which explained an extra amount of variance in SiN
635 performance. These findings may collectively form a neural substrate of individual differences in
636 SiN understanding ability; the variance in SiN perception may be a matter of both primary
637 processes that extract the signal from noise and later speech recognition processes to extract
638 lexical information from speech signals promptly.

639

640 **Acknowledgments**

641 This work was supported by Department of Defense Hearing Restoration and Rehabilitation
642 Program grant awarded to Choi (W81XWH1910637), NIH T32 (5T32DC000040-24) awarded to
643 Schwalje, and NIDCD P50 (DC000242 31) awarded to Griffiths and McMurray. The authors
644 declare no competing financial interests.

645

646 References

- 647 Allopenna, P.D., Magnuson, J.S., Tanenhaus, M.K., 1998. Tracking the Time Course of Spoken
648 Word Recognition Using Eye Movements: Evidence for Continuous Mapping Models. *J*
649 *Mem Lang* 38, 419-439.
- 650 Anderson, S., Kraus, N., 2010. Objective Neural Indices of Speech-in-Noise Perception. *Trends*
651 *in Amplification Trends in Amplification* 14, 73-83.
- 652 Anderson, S., Parbery-Clark, A., White-Schwoch, T., Kraus, N., 2013. Auditory Brainstem
653 Response to Complex Sounds Predicts Self-Reported Speech-in-Noise Performance.
654 *Journal of Speech, Language, and Hearing Research* 56, 31-43.
- 655 Arnal, L.H., Giraud, A.L., 2012. Cortical oscillations and sensory predictions. *Trends Cogn Sci*
656 16, 390-398.
- 657 Ben-David, B.M., Chambers, C.G., Daneman, M., Pichora-Fuller, M.K., Reingold, E.M.,
658 Schneider, B.A., 2011. Effects of aging and noise on real-time spoken word recognition:
659 evidence from eye movements. *J Speech Lang Hear Res* 54, 243-262.
- 660 Bidelman, G.M., Dexter, L., 2015. Bilinguals at the "cocktail party": dissociable neural activity in
661 auditory-linguistic brain regions reveals neurobiological basis for nonnative listeners'
662 speech-in-noise recognition deficits. *Brain Lang* 143, 32-41.
- 663 Bidelman, G.M., Howell, M., 2016. Functional changes in inter- and intra-hemispheric cortical
664 processing underlying degraded speech perception. *Neuroimage* 124, 581-590.
- 665 Brainard, D.H., 1997. The Psychophysics Toolbox. *Spat Vis* 10, 433-436.
- 666 Bregman, A.S., 1999. Auditory scene analysis : the perceptual organization of sound. MIT
667 Press, Cambridge, Mass.
- 668 Bressler, S., Goldberg, H., Shinn-Cunningham, B., 2017. Sensory coding and cognitive
669 processing of sound in Veterans with blast exposure. *Hear Res* 349, 98-110.
- 670 Brouwer, S., Bradlow, A.R., 2016. The Temporal Dynamics of Spoken Word Recognition in
671 Adverse Listening Conditions. *J Psycholinguist Res* 45, 1151-1160.
- 672 Calderone, D.J., Lakatos, P., Butler, P.D., Castellanos, F.X., 2014. Entrainment of neural
673 oscillations as a modifiable substrate of attention. *Trends in Cognitive Sciences Trends in*
674 *Cognitive Sciences* 18, 300-309.
- 675 Cöeponien, R., Cheour, M., Näätänen, R., 1998. Interstimulus interval and auditory event-
676 related potentials in children: evidence for multiple generators. *Electroencephalography and*
677 *Clinical Neurophysiology/Evoked Potentials Section Electroencephalography and Clinical*
678 *Neurophysiology/Evoked Potentials Section* 108, 345-354.
- 679 Choi, I., Kim, S., Choi, I., Schwalje, A., 2020. Cortical dynamics of speech-in-noise
680 understanding. *Acoust. Sci. Technol. Acoustical Science and Technology* 41, 400-403.
- 681 Choi, I., Rajaram, S., Varghese, L.A., Shinn-Cunningham, B.G., 2013. Quantifying attentional
682 modulation of auditory-evoked cortical responses from single-trial electroencephalography.
683 *Front Hum Neurosci* 7, 115.
- 684 Choi, I., Wang, L., Bharadwaj, H., Shinn-Cunningham, B., 2014. Individual differences in
685 attentional modulation of cortical responses correlate with selective attention performance.
686 *Hear Res* 314, 10-19.
- 687 Dahan, D., Gareth Gaskell, M., 2007. The temporal dynamics of ambiguity resolution: Evidence
688 from spoken-word recognition. *J Mem Lang* 57, 483-501.
- 689 Dahan, D., Magnuson, J.S., 2006. Spoken word recognition. *Handbook of psycholinguistics*.
690 Elsevier, pp. 249-283.
- 691 Dai, L., Shinn-Cunningham, B.G., 2016. Contributions of sensory coding and attentional control
692 to individual differences in performance in spatial auditory selective attention tasks. *Front*
693 *Hum Neurosci* 10, 530.
- 694 Dai, L., Shinn-Cunningham, B.G., Best, V., 2018. Sensorineural hearing loss degrades
695 behavioral and physiological measures of human spatial selective auditory attention. *Proc.*

- 696 Natl. Acad. Sci. U. S. A. Proceedings of the National Academy of Sciences of the United
697 States of America 115, E3286-E3295.
- 698 Dale, A.M., Liu, A.K., Fischl, B.R., Buckner, R.L., Belliveau, J.W., Lewine, J.D., Halgren, E.,
699 2000. Dynamic statistical parametric mapping: combining fMRI and MEG for high-resolution
700 imaging of cortical activity. *Neuron* 26, 55-67.
- 701 Darwin, C.J., 1997. Auditory grouping. *Trends Cogn Sci* 1, 327-333.
- 702 Davis, M.H., 2016. The Neurobiology of Lexical Access. *Neurobiology of Language*. Elsevier,
703 pp. 541-555.
- 704 Davis, M.H., Gaskell, M.G., 2009. A complementary systems account of word learning: neural
705 and behavioural evidence. *Philosophical transactions of the Royal Society of London*.
706 Series B, Biological sciences. 364, 3773.
- 707 Davis, M.H., Johnsrude, I.S., 2003. Hierarchical processing in spoken language comprehension.
708 *The Journal of neuroscience : the official journal of the Society for Neuroscience* 23, 3423-
709 3431.
- 710 Ding, N., Melloni, L., Zhang, H., Tian, X., Poeppel, D., 2016. Cortical tracking of hierarchical
711 linguistic structures in connected speech. *Nat Neurosci Nature Neuroscience* 19, 158-164.
- 712 Ding, N., Simon, J.Z., Chatterjee, M., 2014. Robust cortical entrainment to the speech envelope
713 relies on the spectro-temporal fine structure. *NeuroImage NeuroImage* 88, 41-46.
- 714 Du, Y., Buchsbaum, B.R., Grady, C.L., Alain, C., 2014. Noise differentially impacts phoneme
715 representations in the auditory and speech motor systems. *Proc Natl Acad Sci U S A* 111,
716 7126-7131.
- 717 Du, Y., Buchsbaum, B.R., Grady, C.L., Alain, C., 2016. Increased activity in frontal motor cortex
718 compensates impaired speech perception in older adults. *Nat Commun* 7, 12241.
- 719 Etard, O., Reichenbach, T., 2019. Neural Speech Tracking in the Theta and in the Delta
720 Frequency Band Differentially Encode Clarity and Comprehension of Speech in Noise. *The*
721 *Journal of neuroscience : the official journal of the Society for Neuroscience* 39, 5750-5759.
- 722 Farris-Trimble, A., McMurray, B., 2013. Test-retest reliability of eye tracking in the visual world
723 paradigm for the study of real-time spoken word recognition. *Journal of Speech Language*
724 *and Hearing Research* 56, 1328-1345.
- 725 Farris-Trimble, A., McMurray, B., Cigrand, N., Tomblin, J.B., 2014. The process of spoken word
726 recognition in the face of signal degradation. *Journal of Experimental Psychology: Human*
727 *Perception and Performance Journal of Experimental Psychology: Human Perception and*
728 *Performance* 40, 308-327.
- 729 Fischl, B., van der Kouwe, A., Destrieux, C., Halgren, E., Ségonne, F., Salat, D.H., Busa, E.,
730 Seidman, L.J., Goldstein, J., Kennedy, D., Caviness, V., Makris, N., Rosen, B., Dale, A.M.,
731 2004. Automatically parcellating the human cerebral cortex. *Cereb Cortex* 14, 11-22.
- 732 Frey, J.N., Mainy, N., Lachaux, J.P., Muller, N., Bertrand, O., Weisz, N., 2014. Selective
733 Modulation of Auditory Cortical Alpha Activity in an Audiovisual Spatial Attention Task.
734 *JOURNAL OF NEUROSCIENCE* 34, 6634.
- 735 Friston, K., Harrison, L., Daunizeau, J., Kiebel, S., Phillips, C., Trujillo-Barreto, N., Henson, R.,
736 Flandin, G., Mattout, J., 2008. Multiple sparse priors for the M/EEG inverse problem.
737 *Neuroimage* 39, 1104-1120.
- 738 Giraud, A.L., Poeppel, D., 2012. Cortical oscillations and speech processing: emerging
739 computational principles and operations. *Nature neuroscience* 15, 511-517.
- 740 Goldberg, H.R., Choi, I., Varghese, L.A., Bharadwaj, H., Shinn-Cunningham, B.G., 2014.
741 Auditory attention in a dynamic scene: Behavioral and electrophysiological correlates. *The*
742 *Journal of the Acoustical Society of America The Journal of the Acoustical Society of*
743 *America* 135, 2415.
- 744 Golumbic, E.M.Z., Ding, N., Bickel, S., Lakatos, P., Schevon, C.A., McKhann, G.M., Goodman,
745 R.R., Emerson, R., Mehta, A.D., Simon, J.Z., Poeppel, D., Schroeder, C.E., 2013.

- 746 Mechanisms Underlying Selective Neuronal Tracking of Attended Speech at a Cocktail
747 Party. *NEURON* Neuron 77, 980-991.
- 748 Gow, D.W., Jr., 2012. The cortical organization of lexical knowledge: a dual lexicon model of
749 spoken language processing. *Brain Lang* 121, 273-288.
- 750 Gramfort, A., Luessi, M., Larson, E., Engemann, D.A., Strohmeier, D., Brodbeck, C., Goj, R.,
751 Jas, M., Brooks, T., Parkkonen, L., Hamalainen, M., 2013. MEG and EEG data analysis
752 with MNE-Python. *Front Neurosci* 7, 267.
- 753 Gramfort, A., Luessi, M., Larson, E., Engemann, D.A., Strohmeier, D., Brodbeck, C., Parkkonen,
754 L., Hamalainen, M.S., 2014. MNE software for processing MEG and EEG data.
755 *Neuroimage* 86, 446-460.
- 756 Hämäläinen, M.S., Sarvas, J., 1989. Realistic conductivity geometry model of the human head
757 for interpretation of neuromagnetic data. *IEEE Trans Biomed Eng* 36, 165-171.
- 758 Harris, R.W., Swenson, D.W., 1990. Effects of reverberation and noise on speech recognition
759 by adults with various amounts of sensorineural hearing impairment. *Audiology* 29, 314-
760 321.
- 761 Herrmann, B.r., Henry, M.J., Haegens, S., Obleser, J., 2016. Temporal expectations and neural
762 amplitude fluctuations in auditory cortex interactively influence perception. *YNIMG*
763 *NeuroImage* 124, 487-497.
- 764 Hickok, G., Poeppel, D., 2007. The cortical organization of speech processing. *Nat Rev*
765 *Neurosci* 8, 393-402.
- 766 Hillyard, S.A., Hink, R.F., Schwent, V.L., Picton, T.W., 1973. Electrical signs of selective
767 attention in the human brain. *Science* 182, 177-180.
- 768 Hillyard, S.A., Vogel, E.K., Luck, S.J., 1998. Sensory gain control (amplification) as a
769 mechanism of selective attention: electrophysiological and neuroimaging evidence. *Philos*
770 *Trans R Soc Lond B Biol Sci* 353, 1257-1270.
- 771 Holmes, E., Griffiths, T.D., 2019. 'Normal' hearing thresholds and fundamental auditory grouping
772 processes predict difficulties with speech-in-noise perception. *Scientific reports* 9.
- 773 Hornickel, J., Skoe, E., Nicol, T., Zecker, S., Kraus, N., 2009. Subcortical differentiation of stop
774 consonants relates to reading and speech-in-noise perception. *Proc Natl Acad Sci U S A*
775 106, 13022-13027.
- 776 Huettig, F., Altmann, G.T., 2005. Word meaning and the control of eye fixation: semantic
777 competitor effects and the visual world paradigm. *Cognition* 96, B23-32.
- 778 Kislser, T., Reichel, U.D., Sciel, F., 2017. Multilingual processing of speech via web services.
779 *Computer Speech & Language* 45, 326-347.
- 780 Koessler, L., Maillard, L., Benhadid, A., Vignal, J.P., Felblinger, J., Vespignani, H., Braun, M.,
781 2009. Automated cortical projection of EEG sensors: Anatomical correlation via the
782 international 10-10 system. *YNIMG NeuroImage* 46, 64-72.
- 783 Kong, Y.Y., Somarowthu, A., Ding, N., 2015. Effects of Spectral Degradation on Attentional
784 Modulation of Cortical Auditory Responses to Continuous Speech. *JARO -NEW YORK-* 16,
785 783-796.
- 786 Lange, K., 2009. Brain correlates of early auditory processing are attenuated by expectations for
787 time and pitch. *Brain and cognition*. 69, 127.
- 788 Lee, A.K.C., Rajaram, S., Xia, J., Bharadwaj, H., Larson, E., Hämäläinen, M.S., Shinn-
789 Cunningham, B.G., 2013. Auditory Selective Attention Reveals Preparatory Activity in
790 Different Cortical Regions for Selection Based on Source Location and Source Pitch. *Front.*
791 *Neurosci. Frontiers in Neuroscience* 6, 190.
- 792 Liberman, M.C., Epstein, M.J., Cleveland, S.S., Wang, H., Maison, S.F., 2016. Toward a
793 differential diagnosis of hidden hearing loss in humans. *PLoS One* 11, 1-16.
- 794 Luck, S.J., 2014. An introduction to the event-related potential technique, Second edition. ed.
795 The MIT Press, Cambridge, Massachusetts.

- 796 Magnuson, J.S., Dixon, J.A., Tanenhaus, M.K., Aslin, R.N., 2007. The dynamics of lexical
797 competition during spoken word recognition. *Cogn Sci* 31, 133-156.
- 798 Maris, E., Oostenveld, R., 2007. Nonparametric statistical testing of EEG- and MEG-data. *J*
799 *Neurosci Methods* 164, 177-190.
- 800 McMurray, B., Ellis, T.P., Apfelbaum, K.S., 2019. How Do You Deal With Uncertainty? Cochlear
801 Implant Users Differ in the Dynamics of Lexical Processing of Noncanonical Inputs. *Ear and*
802 *hearing* 40, 961-980.
- 803 McMurray, B., Farris-Trimble, A., Rigler, H., 2017. Waiting for lexical access: Cochlear implants
804 or severely degraded input lead listeners to process speech less incrementally. *Cognition*
805 169, 147-164.
- 806 McMurray, B., Samelson, V.M., Lee, S.H., Bruce Tomblin, J., 2010. Individual differences in
807 online spoken word recognition: Implications for SLI. *Cognitive Psychology Cognitive*
808 *Psychology* 60, 1-39.
- 809 McQueen, J.M., Huettig, F., 2012. Changing only the probability that spoken words will be
810 distorted changes how they are recognized. *The Journal of the Acoustical Society of*
811 *America The Journal of the Acoustical Society of America* 131, 509-517.
- 812 Mesgarani, N., Chang, E.F., 2012. Selective cortical representation of attended speaker in multi-
813 talker speech perception. *Nature* 485, 233-236.
- 814 Moore, B.C., 1990. Co-modulation masking release: spectro-temporal pattern analysis in
815 hearing. *British journal of audiology* 24, 131-137.
- 816 Myers, E.B., Blumstein, S.E., Walsh, E., Eliassen, J., 2009. Inferior Frontal Regions Underlie the
817 Perception of Phonetic Category Invariance. *PSCI Psychological Science* 20, 895-903.
- 818 Nabelek, A.K., 1988. Identification of vowels in quiet, noise, and reverberation: relationships
819 with age and hearing loss. *J Acoust Soc Am* 84, 476-484.
- 820 O'Sullivan, J., Mesgarani, N., Smith, E., Schevon, C., McKhann, G.M., Sheth, S.A., Herrero, J.,
821 Mehta, A.D., Smith, E., Sheth, S.A., 2019. Hierarchical Encoding of Attended Auditory
822 Objects in Multi-talker Speech Perception. *NEURON Neuron* 104, 1195-1209.e1193.
- 823 Obleser, J., Erb, J., 2020. Neural Filters for Challenging Listening Situations. *The Cognitive*
824 *Neurosciences*. MIT Press.
- 825 Obleser, J., Kayser, C., 2019. Neural Entrainment and Attentional Selection in the Listening
826 Brain. *Trends Cogn Sci* 23, 913-926.
- 827 Ohlenforst, B., Zekveld, A.A., Lunner, T., Wendt, D., Naylor, G., Wang, Y., Versfeld, N.J.,
828 Kramer, S.E., 2017. Impact of stimulus-related factors and hearing impairment on listening
829 effort as indicated by pupil dilation. *Hear Res* 351, 68-79.
- 830 Owens, E., Schubert, E.D., 1977. Development of the California Consonant Test. *J Speech*
831 *Lang Hear Res.* 20, 463-474.
- 832 Parbery-Clark, A., Skoe, E., Kraus, N., 2009. Musical experience limits the degradative effects
833 of background noise on the neural processing of sound. *J Neurosci* 29, 14100-14107.
- 834 Pelli, D.G., 1997. The VideoToolbox software for visual psychophysics: transforming numbers
835 into movies. *Spat Vis* 10, 437-442.
- 836 Rauschecker, J.P., Scott, S.K., 2009. Maps and streams in the auditory cortex: nonhuman
837 primates illuminate human speech processing. *Nature neuroscience* 12, 718-724.
- 838 Rigler, H., Farris-Trimble, A., Greiner, L., Walker, J., Tomblin, J.B., McMurray, B., 2015. The
839 slow developmental timecourse of real-time spoken word recognition. *Developmental*
840 *psychology* 51, 1690-1703.
- 841 Scott, S.K., Johnsrude, I.S., 2003. The neuroanatomical and functional organization of speech
842 perception. *Trends Neurosci* 26, 100-107.
- 843 Shamma, S., Elhilali, M., Ma, L., Micheyl, C., Oxenham, A.J., Pressnitzer, D., Yin, P., Xu, Y.,
844 2013. Temporal coherence and the streaming of complex sounds. *Advances in*
845 *experimental medicine and biology* 787, 535-543.

- 846 Shinn-Cunningham, B.G., 2020. 14 Brain Mechanisms of Auditory Scene Analysis. *The*
847 *Cognitive Neurosciences*, pp. 159-166.
- 848 Shinn-Cunningham, B.G., Best, V., 2008. Selective attention in normal and impaired hearing.
849 *Trends Amplif* 12, 283-299.
- 850 Song, J.H., Skoe, E., Banai, K., Kraus, N., 2011. Perception of Speech in Noise: Neural
851 Correlates. *Journal of Cognitive Neuroscience* 23, 2268-2279.
- 852 Strait, D.L., Kraus, N., 2011. Can you hear me now? Musical training shapes functional brain
853 networks for selective auditory attention and hearing speech in noise. *Frontiers in*
854 *psychology* 2, 113-113.
- 855 Strauss, D.J., Francis, A.L., 2017. Toward a taxonomic model of attention in effortful listening.
856 *Cognitive, Affective, & Behavioral Neuroscience* 17, 809-825.
- 857 Taylor, J.S.H., Rastle, K., Davis, M.H., 2013. Can Cognitive Models Explain Brain Activation
858 During Word and Pseudoword Reading? A Meta-Analysis of 36 Neuroimaging Studies.
859 *PSYCHOLOGICAL BULLETIN* 139, 766-791.
- 860 Teki, S., Kumar, S., Von Kriegstein, K., Griffiths, T.D., Chait, M., 2011. Brain bases for auditory
861 stimulus-driven figure-ground segregation. *J. Neurosci. Journal of Neuroscience* 31, 164-
862 171.
- 863 Tong, Y., Chen, Q., Nichols, T.E., Rasetti, R., Callicott, J.H., Berman, K.F., Weinberger, D.R.,
864 Mattay, V.S., 2016. Seeking optimal region-of-interest (ROI) single-value summary
865 measures for fMRI studies in imaging genetics. *PLoS One* 11, 1-20.
- 866 Tune, S., Alavash, M., Fiedler, L., Obleser, J., 2020. Neural attention filters do not predict
867 behavioral success in a large cohort of aging listeners. *bioRxiv*, 2020.2005.2020.105874.
- 868 Vaden, K.I., Jr., Kuchinsky, S.E., Ahlstrom, J.B., Dubno, J.R., Eckert, M.A., 2015. Cortical
869 activity predicts which older adults recognize speech in noise and when. *J Neurosci* 35,
870 3929-3937.
- 871 Viswanathan, V., Bharadwaj, H.M., Shinn-Cunningham, B.G., 2019. Electroencephalographic
872 Signatures of the Neural Representation of Speech during Selective Attention. *eNeuro* 6.
- 873 Weber, A., Scharenborg, O., 2012. Models of spoken-word recognition Models of word
874 recognition. *WIREs Cogn Sci Wiley Interdisciplinary Reviews: Cognitive Science* 3, 387-
875 401.
- 876 Wong, P.C., Jin, J.X., Gunasekera, G.M., Abel, R., Lee, E.R., Dhar, S., 2009. Aging and cortical
877 mechanisms of speech perception in noise. *Neuropsychologia* 47, 693-703.
- 878 Wong, P.C., Uppunda, A.K., Parrish, T.B., Dhar, S., 2008. Cortical mechanisms of speech
879 perception in noise. *J Speech Lang Hear Res* 51, 1026-1041.
- 880 Zekveld, A.A., Heslenfeld, D.J., Festen, J.M., Schoonhoven, R., 2006. Top-down and bottom-up
881 processes in speech comprehension. *Neuroimage* 32, 1826-1836.
- 882
- 883

884 **Figure Legends**

885 **Figure 1.** Trial and stimulus structure. Every trial starts with the cue phrase “check the word.” A
886 target word starts 1 second after the noise onset. Four choices are given after the word ends;
887 subjects select the correct answer with a keypad. No feedback is given. The noise level is
888 manipulated to create high (+3 dB) and low (-3 dB) SNR conditions. Subjects complete 50 trials
889 for each condition.
890

891 **Figure 2.** Behavioral results. **A.** Summary of behavioral performance for the two conditions (+3
892 and -3 dB SNR). Boxes denote the 25th – 75th percentile range; the horizontal bars in the center
893 denote the median; the ranges are indicated by vertical dashed lines. Solid lines connect points
894 for the same subject in different conditions. **B.** Average accuracy as a function of reaction time
895 in +3 dB SNR condition. **C.** Average accuracy and reaction time in -3 dB SNR condition.
896

897 **Figure 3.** Sensor-level event-related potential (ERP) with the topographical layout and cortical
898 maps. The time course of the auditory ERP and its envelope, with the standard error of the
899 mean (± 1 SEM) at the peak amplitude (red color: +3 dB SNR, blue color: -3 dB SNR). An
900 asterisk shows a significant difference in the amplitude between +3 and -3 dB SNR conditions
901 (paired *t*-test). Top panels show peak P2 amplitudes of all electrodes in topographical layouts.
902 The *t*-values from paired *t*-tests between two SNR conditions are also shown as topographies.
903

904 **Figure 4.** Region-of-interest (ROI) based source analysis. **A and B.** The time course of the
905 event-related potential (ERP) envelope, with the standard error of the mean (± 1 SEM), obtained
906 at representative voxels for two ROIs in the left hemisphere (“LH”): supramarginal gyrus (SMG),
907 and the pars opercularis and triangularis of the inferior frontal gyrus (IFG), respectively, in each
908 SNR condition (red color: +3 dB SNR, blue color: -3 dB SNR). The ERP envelope time courses
909 in the corresponding regions in the right hemisphere (“RH”) are also shown for visual
910 comparisons. Asterisks show the timing of a significant difference between +3 and -3 dB SNR
911 conditions (cluster-based permutation test, $p = 0.0020$) at the left SMG. **C and D.** Whole-brain
912 maps showing statistical contrasts (*t*-values obtained from post-hoc paired *t*-tests between the
913 two SNR conditions) of source activation at each voxel at the peak timepoint of the grand-
914 average source time course of each ROI.
915

916 **Figure 5.** Top and second panel show histograms of boundaries between phonemes of each
917 stimulus. The third panel shows superimposed temporal envelopes extracted from waveforms of
918 the 100 words.
919

920 **Figure 6.** Individual differences in speech-in-noise processing. **A.** Regression coefficients and
921 their standard errors. **B.** A scatter plot showing the relationship between predicted and
922 measured accuracy in -3 dB SNR condition. **C.** Post-hoc correlation analyses: Raw correlations
923 between each independent variable and the dependent variable. **D.** Left: Relationship between
924 independent variables shows no correlation between internal SNR and evoked source current at
925 the left supramarginal gyrus (SMG). Right: Semi-partial correlation between SMG evoked
926 source current and the residual of accuracy after regressing out internal SNR. **E.** The time
927 course of the auditory event-related potential and its envelope, with the standard error of the
928 mean (± 1 SEM) at the peak magnitude in -3 dB SNR condition (red color: good performers, blue

929 color: poor performers). An asterisk shows a significant difference in the magnitude between
930 two groups (two-sample *t*-test).
931