1    **Cortical Representations of Speech in a Multi-talker Auditory Scene**

2    Krishna C. Puvvada[a], Jonathan Z. Simon[a,b,c]

3    [a] Department of Electrical & Computer Engineering, University of Maryland,

4     College Park, MD 20742

5    [b] Department of Biology, University of Maryland, College Park, MD 20742

6    [c] Institute for Systems Research, University of Maryland, College Park, MD 20742

7
8    Corresponding author: Jonathan Z. Simon (jzsimon@umd.edu), Department of

9    Electrical & Computer Engineering, University of Maryland, 8223 Paint Branch

10   Dr., College Park, MD 20742.

11
12   Abbreviated title: **Foreground and Background Speech Representations**

13   Number of pages: 42

14   Number of figures: 5

15   Number of words in Abstract: 213

16   Number of words in Introduction: 621

17   Number of words in Discussion: 1485

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

22    **Abstract**

23    The ability to parse a complex auditory scene into perceptual objects is facilitated

24    by a hierarchical auditory system. Successive stages in the hierarchy transform an

25    auditory scene of multiple overlapping sources, from peripheral tonotopically-

26    based representations in the auditory nerve, into perceptually distinct auditory-

27    objects based representation in auditory cortex. Here, using magnetoencephalo-

28    graphy (MEG) recordings from human subjects, we investigate how a complex

29    acoustic scene consisting of multiple speech sources is represented in distinct

30    hierarchical stages of auditory cortex. Using systems-theoretic methods of

31    stimulus reconstruction, we show that the primary-like areas in auditory cortex

32    contain dominantly spectro-temporal based representations of the entire auditory

33    scene. Here, both attended and ignored speech streams are represented with almost

34    equal fidelity, and a global representation of the full auditory scene with all its

35    streams is a better candidate neural representation than that of individual streams

36    being represented separately. In contrast, we also show that higher order auditory

37    cortical areas represent the attended stream separately, and with significantly

38    higher fidelity, than unattended streams. Furthermore, the unattended background

39    streams are more faithfully represented as a single unsegregated background

40    object rather than as separated objects. Taken together, these findings demonstrate

41    the progression of the representations and processing of a complex acoustic scene

42    up through the hierarchy of human auditory cortex.

2

43 **Significance Statement:**

44 Using magnetoencephalography (MEG) recordings from human listeners in a

45 simulated cocktail party environment, we investigate how a complex acoustic

46 scene consisting of multiple speech sources is represented in separate hierarchical

47 stages of auditory cortex. We show that the primary-like areas in auditory cortex

48 use a dominantly spectro-temporal based representation of the entire auditory

49 scene, with both attended and ignored speech streams represented with almost

50 equal fidelity. In contrast, we show that higher order auditory cortical areas

51 represent an attended speech stream separately from, and with significantly higher

52 fidelity than, unattended speech streams. Furthermore, the unattended background

53 streams are represented as a single undivided background object rather than as

54 distinct background objects.

55

**Introduction**

56

57 Individual sounds originating from multiple sources in a complex auditory scene

58 mix linearly and irreversibly before they enter the ear, yet are perceived as distinct

59 objects by the listener (Cherry, 1953; Bregman, 1994; McDermott, 2009). The

60 separation, or rather individual re-creation, of such linearly mixed original sound

61 sources is a mathematically ill-posed question, yet the brain nevertheless routinely

62 performs this task with ease. The neural mechanisms by which this perceptual 'un-

63 mixing' of sounds occur, the collective cortical representations of the auditory

64 scene and its constituents, and the role of attention in both, are key problems in

65 contemporary auditory neuroscience.

66 It is known that auditory processing in primate cortex is hierarchical (Davis

67 and Johnsrude, 2003; Hickok and Poeppel, 2007; Rauschecker and Scott, 2009;

68 Okada et al., 2010; Peelle et al., 2010) with subcortical areas projecting onto the

69 core areas of auditory cortex, and from there, on to belt, parabelt and additional

70 auditory areas (Kaas and Hackett, 2000). Sound entering the ear reaches different

71 anatomical/functional areas of auditory cortex with different latencies (Recanzone

72 et al., 2000; Nourski et al., 2014). Due to this serial component of auditory

73 processing, the hierarchy of processing can be described by both anatomy and

74 latency, of which the latter may be exploited using the high temporal fidelity of

75 non-invasive magnetoencephalography (MEG) neural recordings.

4

76      In selective listening experiments using natural speech and MEG, the two

77  major neural responses known to track the speech envelope are the $M50_{TRF}$ and

78  $M100_{TRF}$, with respective latencies of $30 - 80$ ms and $80 - 150$ ms, of which the

79  dominant neural sources are, respectively, Heschl's gyrus (HG) and Planum

80  temporale (PT) (Steinschneider et al., 2011; Ding and Simon, 2012a).

81  Posteromedial HG is the site of core auditory cortex; PT contains both belt and

82  parabelt auditory areas (here collectively referred to as higher-order areas)

83  (Griffiths and Warren, 2002; Sweet et al., 2005). Hence the earlier neural

84  responses are dominated by core auditory cortex, and the later are dominated by

85  higher-order areas. To better understand the neural mechanisms of auditory scene

86  analysis, it is essential to understand how the cortical representations of a complex

87  auditory scene change from the core to the higher order auditory areas.

88      One topic of interest is whether the brain maintains distinct neural

89  representations for each unattended source (in addition to the representation of the

90  attended source), or if all unattended sources are represented collectively as a

91  single monolithic background object. A common paradigm used to investigate the

92  neural mechanisms underlying auditory scene analysis employs a pair of speech

93  streams, of which one is attended, which then leaves the other speech stream

94  remaining as the background (Ding and Simon, 2012a; Mesgarani and Chang,

95  2012; Zion Golumbic et al., 2013b). This results in a limitation, which cannot

5

96    address the question of distinct vs. collective neural representations for unattended

97    sources. This touches on the long-standing debate of whether auditory object

98    segregation is pre-attentive or it is actively influenced by attention (Carlyon, 2004;

99    Sussman et al., 2005; Shinn-Cunningham, 2008; Shamma et al., 2011). Evidence

100   for segregated neural representations of background streams would support the

101   former, whereas a lack of segregated background objects would support the latter.

102          To address these issues, we use MEG to investigate a variety of potential

103   cortical representations of the elements of a multi-talker auditory scene. We test

104   two major hypotheses: that the dominant representation in core auditory cortex is

105   of the physical acoustics, not of separated auditory objects; and that once object-

106   based representations emerge in higher order auditory areas, the unattended

107   contributions to the auditory scene are represented collectively as a single

108   background object. The methodological approach employs the linear systems

109   methods of stimulus prediction and MEG response reconstruction (Ding and

110   Simon, 2012a; Mesgarani and Chang, 2012; Di Liberto et al., 2015).

111

112   **Materials & Methods**:

113   ***Subjects & Experimental Design*** Nine normal-hearing, young adults (6 Female)

114   participated in the experiment. All subjects were paid for their participation. The

115   experimental procedures were approved by the University of Maryland

6

116     Institutional Review Board. Subjects listened to a mixture of three speech

117     segments spoken by, respectively, a male adult, female adult and a child speaker.

118     The three speech segments were mixed into a single audio channel with equal

119     perceptual loudness. All three speech segments were taken from public domain

120     narration of Grimms' Fairy Tales by Jacob & Wilhelm Grimm

121     (https://librivox.org/fairy-tales-by-the-brothers-grimm/). Periods of silence longer

122     than 300 ms were replaced by a shorter gap whose duration was chosen randomly

123     between 200 ms and 300 ms. The audio signal was low-pass filtered below 4 kHz.

124     In first of three conditions, the subjects were asked to attend to the child speaker,

125     while ignoring the other two (i.e., child speaker as target, with male and female

126     adult speakers as background). In condition two, during which the same mixture

127     was played as in condition one, the subjects were instead asked to attend to the

128     male adult speaker (with female adult and child speakers as background).

129     Similarly, in condition three, the target was switched to the female adult speaker.

130     Each condition was repeated three times successively, producing three trials per

131     condition. The presentation order of the three conditions was counterbalanced

132     across subjects. Each trial was of 220 s duration, divided into two 110 s sections,

133     to reduce listener fatigue. To help participants attend to the correct speaker, the

134     first 30 s of each section was replaced by the clean recording of the target speaker

135     alone, followed by a 5 s upward linear ramp of the background speakers.

7

136    Recordings of this first 35 s of each segment were not included in any analysis. To

137    further encourage the subjects to attend to the correct speaker, a target-word was

138    set before each trial and the subjects were asked to count the number of

139    occurrences of the target-word in the speech of the attended speaker. Additionally,

140    after each condition, the subject was asked to recount a short summary of the

141    attended narrative. The subjects were required to close their eyes while listening.

142    Before the main experiment, 100 repetitions of a 500-Hz tone pip were presented

143    to each subject to elicit the M100 response, a reliable auditory response occurring

144    ~100 ms after the onset of a tone pip. This data was used check whether any

145    potential subjects gave abnormal auditory responses, but no subjects were

146    excluded based on this criterion.

147

148    ***Data recording and pre-processing*** MEG recordings were conducted using a 160-

149    channel whole-head system (Kanazawa Institute of Technology, Kanazawa,

150    Japan). Its detection coils are arranged in a uniform array on a helmet-shaped

151    surface of the bottom of the dewar, with ~25 mm between the centers of two

152    adjacent 15.5-mm-diameter coils. Sensors are configured as first-order axial

153    gradiometers with a baseline of 50 mm; their field sensitivities are 5 fT/$\sqrt{\ }$Hz or

154    better in the white noise region. Subjects lay horizontally in a dimly lit

155    magnetically shielded room (Yokogawa Electric Corporation). Responses were

8

156 recorded with a sampling rate of 1 kHz with an online 200-Hz low-pass filter and

157 60 Hz notch filter. Three reference magnetic sensors and three vibrational sensors

158 were used to measure the environmental magnetic field and vibrations. The

159 reference sensor recordings were utilized to reduce environmental noise from the

160 MEG recordings using the Time-Shift PCA method (de Cheveigne and Simon,

161 2007). Additionally, MEG recordings were decomposed into virtual sensors/

162 components using denoising source separation (DSS) (Särelä and Valpola, 2005;

163 de Cheveigne and Simon, 2008; de Cheveigne and Parra, 2014), a blind source

164 separation method that enhances neural activity consistent over trials. Specifically,

165 DSS decomposes the multichannel MEG recording into temporally uncorrelated

166 components, where each component is determined by maximizing its trial-to-trial

167 reliability, measured by the correlation between the responses to the same stimulus

168 in different trials. To reduce the computational complexity, for all further analysis

169 the 157 MEG sensors were reduced, using DSS, to 4 components in each

170 hemisphere. Also, both stimulus envelope and MEG responses were band pass

171 filtered between 1 – 8 Hz (delta and theta bands), which correspond to the slow

172 temporal modulations in speech (Ding and Simon, 2012b, a).

173

174 ***Terminology and Notation*** As specified in the stimulus description, in each

175 condition the subject attends to one among the three speech streams. The envelope

9

176    of attended speech stream is referred to as the 'foreground' and the envelope of

177    each of the two unattended speech streams is referred to as the 'individual

178    background'. In contrast, the envelope of the entire unattended part of the

179    stimulus, comprising *both* unattended speech streams, is referred to as the

180    'combined background'. The envelope of entire acoustic stimulus or auditory

181    scene, comprising of all the three speech streams is referred to as the 'acoustic

182    scene'. Thus, if $S_a, S_b, S_c$ are three speech stimuli, $Env(S_a + S_b + S_c)$ is the

183    acoustic scene. In contrast, the sum of envelopes of three speech streams,

184    $Env(S_a) + Env(S_b) + Env(S_c)$, is referred to as the 'sum of streams', and the

185    two are not mathematically equal: even though both are functions of the same

186    stimuli, they differ due to the non-linear nature of a signal envelope (the linear

187    correlation between the acoustic scene and the sum of streams is typically ~0.75).

188        Neural responses with latencies less than ~85 ms (typically originating

189    from core auditory areas) are referred to here as 'early neural responses' and

190    responses with latencies more than ~85 ms (typically from higher-order auditory

191    areas) (Ahveninen et al., 2011; Okamoto et al., 2011; Steinschneider et al., 2011)

192    are referred to as 'late neural responses'.

193

194    ***Temporal Response Function*** In an auditory scene with a single talker, the

195    relation between MEG neural response and the presented speech stimuli can be

196    modeled using a linear temporal response function (TRF) as

$$r(t) = \sum_{\tau} s(t - \tau)TRF(\tau) + \varepsilon(t) \tag{1}$$

197    where $t = 0,1, \dots, T$ is time, $r(t)$ is the response from any individual sensor or

198    DSS component, $s(t)$ is the stimulus envelope in decibels, $TRF(t)$ is the TRF

199    itself, and $\epsilon(t)$ is residual response waveform not explained by the TRF model

200    (Ding and Simon, 2012b). The envelope is extracted by averaging the auditory

201    spectrogram, (Chi et al., 2005) along the spectral dimension. The TRF is estimated

202    using boosting with 10-fold cross-validation (David et al., 2007). In case of single

203    speech stimuli, the TRF is typically characterized by a positive peak between 30

204    ms and 80 ms and a negative peak between 90 ms and 130 ms, referred to as

205    M50$_{\text{TRF}}$ and M100$_{\text{TRF}}$ respectively (Ding and Simon, 2012a) (positivity/negativity

206    of the magnetic field is by convention defined to agree with the corresponding

207    electroencephalography[EEG] peaks). Success/accuracy of the linear model is

208    evaluated by how well it predicts neural responses, as measured by the proportion

209    of the variance explained: the square of the Pearson correlation coefficient

210    between the MEG measurement and the TRF model prediction.

211         In the case of more than one speaker, the MEG neural response, $r(t)$ can be

212    modeled as the sum of the responses to the individual acoustic sources (Ding and

11

213    Simon, 2012a; Zion Golumbic et al., 2013b), referred to here as the 'Summation

214    model'. For example, with two speech streams, the neural response would be

215    modeled as

$$r(t) = \sum_{\tau} S_a(t - \tau)TRF_a(\tau) + \sum_{\tau} S_b(t - \tau)TRF_b(\tau) + \varepsilon(t) \qquad (2)$$

216

217    where $S_a(t)$ and $S_b(t)$ are the envelopes of the two speech streams, and $TRF_a(t)$,

218    and $TRF_b(t)$ are the TRFs corresponding to each stream. The summation model is

219    easily extended to the case of more than two speech streams, by adding new terms

220    with each new individual speech stream envelope and the corresponding TRF.

221            In addition to the existing summation model, we propose a new encoding-

222    model referred to as the 'Early-late model', which allows one to incorporate the

223    hypothesis that the early neural responses typically represent the entire acoustic

224    scene, but that the later neural responses differentially represent the separated

225    foreground and background.

$$r(t) = \sum_{\tau=0}^{\tau=\tau_1} S_A(t - \tau)TRF_A(\tau) + \sum_{\tau=\tau_1}^{\tau=\tau_2} S_F(t - \tau)TRF_F(\tau) + \sum_{\tau=\tau_1}^{\tau=\tau_2} S_B(t - \tau)TRF_B(\tau) + \epsilon(t) \qquad (3)$$

226

227    where $S_A(t)$ is the (entire) acoustic scene, $S_F(t)$ is the envelope of attended

228    (foreground) speech stream, and $S_B(t)$ is the combined background (i.e., envelope

229    of everything other than attended speech stream in the auditory scene), and

230    $TRF_A(t), TRF_F(t),$ and $TRF_B(t)$ are the corresponding $TRFs$. $\tau_1, \tau_2$ represent the

12

231    boundary values of the integration windows for early and late neural responses

232    respectively.

233          The explanatory power of different models, such as the Summation and

234    Early-late models, can be ranked by comparing the accuracy of their response

235    predictions (illustrated in Figure 1, left).

236

237                               (Figure 1 about here)

238

239    ***Decoding speech from neural responses*** While the TRF/encoding analysis

240    described in the previous section predicts neural response from the stimulus,

241    decoding analysis reconstructs the stimulus based on the neural response. Thus,

242    decoding analysis complements the TRF analysis (Mesgarani et al., 2009).

243    Mathematically the envelope reconstruction/decoding operation can be formulated

244    as

$$E(t) = \sum_{k=1}^{N} \sum_{\tau=\tau_b}^{\tau_e} M_k(t+\tau)D_k(\tau) + \epsilon(t) \qquad (4)$$

245

246    where $E(t)$ is the reconstructed envelope, $M_k(t)$ is the MEG recording (neural

247    response) from sensor/component $k$, and $D_k(t)$ is the linear decoder for

248    sensor/component $k$. The times $\tau_b$ and $\tau_e$ denote the beginning and end times of

249    the integration window. By appropriately choosing the values of $\tau_b$ and $\tau_e$,

13

250    envelope reconstructions using neural responses from any desired time window

251    can be compared. The decoder is estimated using boosting analogously to the TRF

252    estimation in the previous section. In the single talker case the envelope is of that

253    talker's speech. In a multi-talker case, the envelope to be reconstructed might be

254    the envelope of the speech of attended talker, or one of the background talkers, or

255    of a mixture of any two or all three talkers, depending on the model under

256    consideration. Chance-level reconstruction (i.e., the noise floor) from a particular

257    neural response is estimated by reconstructing an unrelated stimulus envelope

258    from that neural response. Figure 2 illustrates the distinction between

259    reconstruction of stimulus envelope from early and late responses. The stimulus

260    envelope at time point $t$ can be reconstructed using neural responses from the

261    dashed (early response) window or dotted (late response) window. (While it is true

262    that the late responses to the stimulus at time point $t - \Delta t$ overlap with early

263    responses to the stimulus at time point $t$, the decoder used to reconstruct the

264    stimulus at time point $t$ from early responses is only minimally affected by late

265    responses to the stimulus at time point $t - \Delta t$ when the decoder is estimated by

266    averaging over a long enough duration, e.g., tens of seconds). The cut-off time

267    between early and late responses, $\tau_{boundary}$, was chosen to minimize the overlap

268    between the $M50_{TRF}$ and $M100_{TRF}$ peaks, on a per subject basis, with a typical

269    value being 85 ms. When decoding from early responses only, the time window of

14

270    integration is from $\tau_b = 0$ to $\tau_e = \tau_{boundary}$. When decoding from late neural

271    responses only, the time window of integration is from $\tau_b = \tau_{boundary}$ to $\tau_e =$

272    500 ms.

273

274                              (Figure 2 about here)

275

276         The robustness of different representations, such as of Foreground vs.

277    Background, can be compared by examining the accuracy of their respective

278    stimulus envelope reconstructions (illustrated in Figure 1, right).

279

280    ***Statistics*** All statistical comparisons reported here are two-tailed permutation tests

281    with $N$=1,000,000 random permutations (within subject). Due to the value of N

282    selected, the smallest accurate $p$ value that can be reported is $2 \times 1/N$ ($= 2 \times 10^{-6}$; the

283    factor of 2 arises from the two-tailed test) and any $p$ value smaller than $2/N$ is

284    reported as $p < 2 \times 10^{-6}$. The statistical comparison between foreground and

285    individual backgrounds requires special mention, since each listening condition

286    has one foreground but two individual backgrounds. From the perspective of both

287    behavior and task, both the individual backgrounds are interchangeable. Hence,

288    when comparing reconstruction accuracy of foreground vs. individual background

289    the average reconstruction accuracy of the two individual backgrounds is used.

15

290     Finally, Bayes factor analysis is used, when appropriate, to evaluate evidence in

291     favor of null hypothesis, since conventional hypothesis testing is not suitable for

292     such purposes. Briefly, Bayes factor analysis calculates the *posterior odds* i.e., the

293     ratio of $P(H_0|observations)$ to $P(H_1|observations)$, where $H_0$ and $H_1$ are the null

294     and alternate hypotheses respectively.

$$\frac{P(H_0|observations)}{P(H_1|observations)} = \frac{P(observations|H_0)}{P(observations|H_1)} \times \frac{P(H_0)}{P(H_1)} \qquad (5)$$

$$= BF_{01} \times \frac{P(H_0)}{P(H_1)} \qquad (6)$$

295     The ratio of $P(observations|H_0)$ and $P(observations|H_1)$ is denoted as the Bayes

296     factor, $BF_{01}$. Then, under the assumption of equal priors ($P(H_0) = P(H_1)$), the

297     posterior odds reduces to $BF_{01}$. A $BF_{01}$ value of 10 indicates that the data is ten

298     times more likely to occur under the null hypothesis than the alternate hypothesis;

299     conversely, a $BF_{01}$ value of 0.1 indicates that the data is 10 times more likely to

300     occur under the alternate hypothesis than the null hypothesis. Conventionally, a

301     $BF_{01}$ value between 3 and 10 is considered as moderate evidence in favor of the

302     null hypothesis, and a value between 10 and 30 is considered strong evidence;

303     conversely, a $BF_{01}$ value between 1/3 & 1/10  (respectively 1/10 & 1/30) is

304     considered moderate (respectively strong) evidence for the alternate hypothesis

305     (for more details we refer the reader to Rouder et al. (2009)).

16

306

## Results

### *Stimulus reconstruction from early neural responses*

To investigate the neural representations of the attended vs. unattended speech

streams associated with early auditory areas, i.e., from core auditory cortex,

(Nourski et al., 2014), the temporal envelope of attended (foreground) and

unattended speech streams (individual backgrounds) were reconstructed using

decoders optimized individually for each speech stream. All reconstructions

performed significantly better than chance level (foreground vs. noise, $p < 2 \times 10^{-6}$;

individual background vs. noise, $p < 2 \times 10^{-6}$), indicating that all three speech

streams are represented in early auditory cortex. Figure 3A shows reconstruction

accuracy for foreground vs. individual backgrounds. A permutation test shows no

significant difference between foreground and individual background ($p = 0.21$),

indicating that there is no evidence of significant neural bias for the attended

speech stream over the ignored speech stream, in early neural responses. In fact,

Bayes Factor analysis ($BF_{01} = 4.2$) indicates moderate support in favor of the null

hypothesis (Rouder et al., 2009), that early neural responses do not distinguish

significantly between attended and ignored speech streams.

324

325                                    (Figure 3 about here)

17

326

327     To test the hypothesis that early auditory areas represent the auditory scene

328     in terms of acoustics, rather than as individual auditory objects, we reconstructed

329     the acoustic scene (the envelope of the sum of all three speech streams) and

330     compared it against the reconstruction of the sum of streams (sum of

331     reconstruction envelopes of each of the three individual speech streams). Separate

332     decoders optimized individually were used to reconstruct the acoustic scene and

333     the sum of streams. As can be seen in Figure 3B, the result shows that the acoustic

334     scene is better reconstructed than the sum of streams ($p < 2 \times 10^{-6}$). This indicates

335     that early auditory cortex is better described as processing the entire acoustic scene

336     rather than processing the separate elements of the scene individually.

337

338     ***Stimulus reconstruction from late neural responses***

339     While the preceding results were based on early cortical processing, the following

340     results are based on late auditory cortical processing (responses with latencies

341     more than ~85 ms). Figure 4A shows the scatter plot of reconstruction accuracy

342     for the foreground vs. individual background envelopes based on late responses. A

343     paired permutation test shows that reconstruction accuracy for the foreground is

344     significantly higher than the background ($p < 2 \times 10^{-6}$). Even though the individual

18

345    backgrounds are not as reliably reconstructed as foreground, their reconstructions

346    are nonetheless significantly better than chance level ($p < 2\times10^{-6}$).

347         In order to distinguish among possible neural representations of the

348    background streams, we compared the reconstructability of the envelope of the

349    entire background as a whole, with the reconstructability of the sum of the

350    envelopes of the (two) backgrounds. If the background is represented as a single

351    auditory object (i.e., "the background"), the reconstruction of the envelope of the

352    entire background should be more faithful than the sum of envelopes of individual

353    backgrounds. In contrast, if the background is represented as distinct auditory

354    objects, each distinguished by its own envelope, the reconstruction of the sum of

355    envelopes of the individual backgrounds should be more faithful. Figure 4B shows

356    the scatter plot of reconstruction accuracy for the envelope of combined

357    background vs. the sum of the envelopes of the individual background streams.

358    Analysis shows that the envelope of the combined background is significantly

359    better represented than the sum of the individual envelopes of the individual

360    backgrounds ($p = 0.012$). As noted previously, the envelope of the combined

361    background is actually strongly correlated with the sum of the envelopes of the

362    individual backgrounds, meaning that finding a significant difference in their

363    reconstruction accuracy is *a priori* unlikely, providing even more credence to the

364    result.

19

365

366                              (Figure 4 about here)

367

368     ***Encoding analysis***

369     Results above from envelope reconstruction suggest that while early neural

370     responses represent the auditory scene in terms of the acoustics, the later neural

371     responses represent the auditory scene in terms of a separated foreground and a

372     single background stream. In order to further test this hypothesis, we use TRF-

373     based encoding analysis to directly compare two different models of auditory

374     scene representations. The two models compared are the standard Summation

375     model (based on parallel representations of all speech streams; see Equation 2) and

376     the new Early-late model (based on an early representation of the entire acoustic

377     scene and late representations of separated foreground and background; see

378     Equation 3). Figure 5 shows the response prediction accuracies for the two

379     models. A permutation test shows that the accuracy of the Early-late model is

380     considerably higher than that of the Summation model ($p < 2\times10^{-6}$). This indicates

381     that a model in which early/core auditory cortex processes the entire acoustic

382     scene but later/higher-order auditory cortex processes the foreground and

383     background separately has more support than the previously employed model of

384     parallel processing of separate streams throughout auditory cortex.


        20

385

386                          (Figure 5 about here)


387    **Discussion**

388    In this study, we used cortical tracking of continuous speech, in a multi-talker

389    scenario, to investigate the neural representations of an auditory scene. Differing

390    latencies of the neural sources processing the same stimuli allow us to separate the

391    source activity temporally, thus enabling the tracking of differing neural

392    representations of the auditory scene. From MEG recordings of subjects

393    selectively attending to one of the three co-located speech streams, we observed

394    that 1) The early neural responses (with short latencies), which originate primarily

395    from core auditory cortex, represent the foreground (attended) and background

396    (ignored) speech streams without any significant difference, whereas the late

397    neural responses (with longer latencies), which originate primarily from higher-

398    order areas of auditory cortex, represent the foreground with significantly higher

399    fidelity than the background; 2) Early neural responses are not only balanced in

400    how they represent the constituent speech streams, but in fact represent the entire

401    acoustic scene holistically, rather than as separately contributing individual

402    perceptual objects; 3) Even though there are two physical speech streams in the

403    background, no neural segregation is observed for the background speech streams.

21

404        It is well established that auditory processing in cortex is performed in a

405    hierarchical fashion, in which an auditory stimulus is processed by different

406    anatomical areas at different latencies (Inui et al., 2006; Nourski et al., 2014).

407    Using this idea to inform the neural decoding/encoding analysis allows the

408    effective isolation of neural signals from a particular cortical area, and thereby the

409    ability to track changes in neural representations as the stimulus processing

410    proceeds along the auditory hierarchy. This time-constrained

411    reconstruction/prediction approach may prove especially fruitful in high-time-

412    resolution/low-spatial-resolution imaging techniques such as MEG and EEG. Even

413    though different response components are generated by different neural sources,

414    standard neural source localization algorithms may perform poorly when different

415    sources are strongly correlated in their responses (Lutkenhoner and Mosher,

416    2007). While the proposed method is not to be viewed as an alternative to source

417    localization methods, it can nonetheless be used to tease apart different

418    components of MEG/EEG response, without explicit source localization.

419        The envelope reconstruction using the early, auditory core, neural response

420    component showed no significant difference between foreground and background,

421    in contrast to reconstruction using the late, higher-order auditory, neural

422    responses, where the foreground is substantially better represented than any

423    individual background. This *decoding* result is in agreement with the *encoding*

22

424    result of (Ding and Simon, 2012a) where the authors showed that the early $M50_{TRF}$

425    component of the temporal response function is not significantly modulated by

426    attention, whereas the late $M100_{TRF}$ component is modulated by attention.

427         Even though there is no significant difference between the ability to

428    reconstruct the foreground and background from early neural responses,

429    nonetheless we observe a non-significant tendency towards an enhanced

430    representation of the foreground (foreground > background, $p = 0.21$). This could

431    be due to task-related plasticity of spectro-temporal receptive fields of neurons in

432    mammalian primary auditory cortex (Fritz et al., 2003), where the receptive fields

433    of neurons are tuned to match the stimulus characteristics of attended sounds. It

434    could also be explained by entrainment (Schroeder and Lakatos, 2009; Zion

435    Golumbic et al., 2012), which postulates that the high excitability periods of

436    neurons become aligned with temporal structure of foreground, thereby enhancing

437    its neural representation.

438         The increase in fidelity of the foreground as the response latency increases,

439    from early neural responses (from core auditory cortex) to late neural responses

440    (from higher-order auditory cortex), indicates a temporal as well as functional

441    hierarchy in cortical processing of auditory scene, from core to higher-order areas

442    in auditory cortex. Similar preferential representation for the attended speech

443    stream has been demonstrated, albeit with only two speech streams, using delta

23

444    and theta band neural responses (Ding and Simon, 2012a; Zion Golumbic et al.,

445    2013a; Zion Golumbic et al., 2013b) as well as high-gamma neural responses

446    (Mesgarani and Chang, 2012; Zion Golumbic et al., 2013a), and using monaural

447    (Ding and Simon, 2012a; Mesgarani and Chang, 2012) as well as audio-visual

448    speech (Zion Golumbic et al., 2013a; Zion Golumbic et al., 2013b).

449        While some researchers suggest a selective entrainment model (Schroeder

450    and Lakatos, 2009; Zion Golumbic et al., 2013b) as the mechanism underlying the

451    selective tracking of attended speech, others suggest a temporal coherence model

452    (Shamma et al., 2011; Ding and Simon, 2012a) as the neuronal mechanism

453    underlying selective tracking. Natural speech is quasi-rhythmic with different

454    dominant rates at syllabic, word and prosodic frequencies. The selective

455    entrainment model suggests that attention causes endogenous low frequency

456    neural oscillations to align with the temporal structure of the attended speech

457    stream, thus aligning the high excitability phases of oscillations with events in

458    attended stream. This effectively forms a mask that favors the attended speech.

459    The temporal coherence model suggests that selective tracking of attended speech

460    is achieved through two stages. First is a cortical filtering stage, where feature

461    selective neurons filter the stimulus producing a multidimensional representation

462    of auditory scene along different feature axes. This is followed by a second stage,

463    coherence analysis, which combines different features streams based on their

24

464    temporal similarity, giving rise to separate perceptions of attended and ignored

465    streams.

466          The representation of an auditory scene in core auditory cortex is here

467    shown to be more spectro-temporal- or acoustic-based than object-based, as

468    demonstrated by the result that the envelope of the auditory scene is better

469    reconstructed than the sum of envelopes of the individual speech streams (e.g.,

470    Figure 3B). This is further supported by the result that the Early-late model

471    predicts MEG neural responses significantly better than Summation model (e.g.,

472    Figure 5). This is consistent with previous studies that demonstrated that neural

473    activity in core auditory cortex was highly sensitive to acoustic characteristics of

474    speech and primarily reflects spectro-temporal attributes of sound (Nourski et al.,

475    2009; Okada et al., 2010; Steinschneider et al., 2014). All these results suggest that

476    early neural responses, primarily from core auditory cortex, reflect an acoustic-

477    based representation rather than object-based. In contrast, Nelken and Bar-Yosef

478    (2008) suggest that neural auditory objects may form as early as primary auditory

479    cortex, and Fritz et al. (2003) show that representations of dynamic sounds in

480    primary auditory cortex are influence by task. It is possible that less complex

481    stimuli are resolved earlier in the hierarchy of auditory pathway (e.g., sounds that

482    can be separated via tonotopy) whereas speech streams, which overlap both

483    spectrally and temporally, are resolved only much later in auditory pathway.

25

484        It is widely accepted that an auditory scene is *perceived* in terms of

485        auditory objects (Bregman, 1994; Griffiths and Warren, 2004; Shinn-Cunningham,

486        2008; Shamma et al., 2011). Ding and Simon (2012b) demonstrated evidence for

487        an object-based cortical representation of an auditory scene, but did not distinguish

488        between early and late neural responses. This, coupled with the result here that

489        early neural responses provide an acoustic, not object-based, representation,

490        strongly suggest that the object-based representation emerges only in the late

491        neural responses/higher-order (belt and parabelt) auditory areas. This is further

492        supported by the observation that acoustic invariance, a property of object-based

493        representation, is observed in higher order areas but not in core auditory cortex

494        (Chang et al., 2010; Okada et al., 2010).

495        When the foreground is represented as an auditory object in late neural

496        responses, the finding that the combined background is better reconstructed than

497        the sum of envelopes of individual backgrounds (Figure 4B) suggests that in late

498        neural responses the background is not represented as separated and distinct

499        auditory objects. This result is consistent with that of Sussman et al. (2005), who

500        reported an unsegregated background when subjects attended to one of three tone

501        streams in the auditory scene. This unsegregated background may be a result of an

502        'analysis-by-synthesis' (Yuille and Kersten, 2006; Poeppel et al., 2008)

503        mechanism, wherein the auditory scene is first decomposed into basic acoustic

26

504    elements, followed by top-down processes that guide the synthesis of the relevant

505    components into a single stream, which then becomes the object of attention. The

506    remainder of the auditory scene would be the unsegregated background, which

507    itself might have the properties of an auditory object. When attention shifts, new

508    auditory objects are correspondingly formed, with the old ones now contributing

509    to the unstructured background. Shamma et al. (2011) suggest that this top down

510    influence acts through the principle of temporal coherence. Between the two

511    opposing views, that streams are formed pre-attentively and that multiple streams

512    can co-exist simultaneously, or that attention is required to form a stream and only

513    that single stream is ever present as separated perceptual entity, these findings lend

514    support to the latter.

515        In summary, these results provide evidence that, in a complex auditory

516    scene with multiple overlapping spectral and temporal sources, the core areas of

517    auditory cortex maintains an acoustic representation of the auditory scene with no

518    significant preference to attended over ignored source, and with no separation into

519    distinct sources. It is only the higher-order auditory areas that provide an object

520    based representation for the foreground, but even there the background remains

521    unsegregated.

522

523

27

524

525 **References**

526 Ahveninen J, Hamalainen M, Jaaskelainen IP, Ahlfors SP, Huang S, Lin FH, Raij T,

527      Sams M, Vasios CE, Belliveau JW (2011) Attention-driven auditory cortex short-

528      term plasticity helps segregate relevant sounds from noise. Proc Natl Acad Sci U S

529      A 108:4182-4187.

530 Bregman AS (1994) Auditory scene analysis: The perceptual organization of sound: MIT

531      press.

532 Carlyon RP (2004) How the brain separates sounds. Trends Cogn Sci 8:465-471.

533 Chang EF, Rieger JW, Johnson K, Berger MS, Barbaro NM, Knight RT (2010)

534      Categorical speech representation in human superior temporal gyrus. Nat Neurosci

535      13:1428-1432.

536 Cherry EC (1953) Some Experiments on the Recognition of Speech, with One and with 2

537      Ears. Journal of the Acoustical Society of America 25:975-979.

538 Chi T, Ru P, Shamma SA (2005) Multiresolution spectrotemporal analysis of complex

539      sounds. J Acoust Soc Am 118:887-906.

540 David SV, Mesgarani N, Shamma SA (2007) Estimating sparse spectro-temporal

541      receptive fields with natural stimuli. Network 18:191-212.

28

542  Davis MH, Johnsrude IS (2003) Hierarchical processing in spoken language

543      comprehension. J Neurosci 23:3423-3431.

544  de Cheveigne A, Simon JZ (2007) Denoising based on time-shift PCA. J Neurosci

545      Methods 165:297-305.

546  de Cheveigne A, Simon JZ (2008) Denoising based on spatial filtering. J Neurosci

547      Methods 171:331-339.

548  de Cheveigne A, Parra LC (2014) Joint decorrelation, a versatile tool for multichannel

549      data analysis. Neuroimage 98:487-505.

550  Di Liberto GM, O'Sullivan JA, Lalor EC (2015) Low-Frequency Cortical Entrainment to

551      Speech Reflects Phoneme-Level Processing. Curr Biol 25:2457-2465.

552  Ding N, Simon JZ (2012a) Emergence of neural encoding of auditory objects while

553      listening to competing speakers. Proc Natl Acad Sci U S A 109:11854-11859.

554  Ding N, Simon JZ (2012b) Neural coding of continuous speech in auditory cortex during

555      monaural and dichotic listening. J Neurophysiol 107:78-89.

556  Fritz J, Shamma S, Elhilali M, Klein D (2003) Rapid task-related plasticity of

557      spectrotemporal receptive fields in primary auditory cortex. Nat Neurosci 6:1216-

558      1223.

559  Griffiths TD, Warren JD (2002) The planum temporale as a computational hub. Trends

560      Neurosci 25:348-353.

29

561    Griffiths TD, Warren JD (2004) What is an auditory object? Nature Reviews

562        Neuroscience 5:887-892.

563    Hickok G, Poeppel D (2007) The cortical organization of speech processing. Nat Rev

564        Neurosci 8:393-402.

565    Inui K, Okamoto H, Miki K, Gunji A, Kakigi R (2006) Serial and parallel processing in

566        the human auditory cortex: a magnetoencephalographic study. Cereb Cortex 16:18-

567        30.

568    Kaas JH, Hackett TA (2000) Subdivisions of auditory cortex and processing streams in

569        primates. Proc Natl Acad Sci U S A 97:11793-11799.

570    Lutkenhoner B, Mosher JC (2007) Source Analysis of Auditory Evoked Potentials and

571        Fields. In: Auditory evoked potentials : basic principles and clinical application

572        (Burkard RF, Eggermont JJ, Don M, eds), pp xix, 731 p., 716 p. of plates.

573        Philadelphia: Lippincott Williams & Wilkins.

574    McDermott JH (2009) The cocktail party problem. Curr Biol 19:R1024-1027.

575    Mesgarani N, Chang EF (2012) Selective cortical representation of attended speaker in

576        multi-talker speech perception. Nature 485:233-236.

577    Mesgarani N, David SV, Fritz JB, Shamma SA (2009) Influence of context and behavior

578        on stimulus reconstruction from neural activity in primary auditory cortex. J

579        Neurophysiol 102:3329-3339.

30

580    Nelken I, Bar-Yosef O (2008) Neurons and objects: the case of auditory cortex. Front

581        Neurosci 2:107-113.

582    Nourski KV, Steinschneider M, McMurray B, Kovach CK, Oya H, Kawasaki H, Howard

583        MA, 3rd (2014) Functional organization of human auditory cortex: investigation of

584        response latencies through direct recordings. Neuroimage 101:598-609.

585    Nourski KV, Reale RA, Oya H, Kawasaki H, Kovach CK, Chen H, Howard MA, 3rd,

586        Brugge JF (2009) Temporal envelope of time-compressed speech represented in the

587        human auditory cortex. J Neurosci 29:15564-15574.

588    Okada K, Rong F, Venezia J, Matchin W, Hsieh IH, Saberi K, Serences JT, Hickok G

589        (2010) Hierarchical organization of human auditory cortex: evidence from acoustic

590        invariance in the response to intelligible speech. Cereb Cortex 20:2486-2495.

591    Okamoto H, Stracke H, Bermudez P, Pantev C (2011) Sound processing hierarchy within

592        human auditory cortex. Journal of Cognitive Neuroscience 23:1855-1863.

593    Peelle JE, Johnsrude IS, Davis MH (2010) Hierarchical processing for speech in human

594        auditory cortex and beyond. Front Hum Neurosci 4:51.

595    Poeppel D, Idsardi WJ, van Wassenhove V (2008) Speech perception at the interface of

596        neurobiology and linguistics. Philos Trans R Soc Lond B Biol Sci 363:1071-1086.

597    Rauschecker JP, Scott SK (2009) Maps and streams in the auditory cortex: nonhuman

598        primates illuminate human speech processing. Nat Neurosci 12:718-724.

31

599    Recanzone GH, Guard DC, Phan ML (2000) Frequency and intensity response properties

600        of single neurons in the auditory cortex of the behaving macaque monkey. J

601        Neurophysiol 83:2315-2331.

602    Rouder JN, Speckman PL, Sun D, Morey RD, Iverson G (2009) Bayesian t tests for

603        accepting and rejecting the null hypothesis. Psychon Bull Rev 16:225-237.

604    Särelä J, Valpola H (2005) Denoising source separation. Journal of Machine Learning

605        Research 6:233-272.

606    Schroeder CE, Lakatos P (2009) Low-frequency neuronal oscillations as instruments of

607        sensory selection. Trends Neurosci 32:9-18.

608    Shamma SA, Elhilali M, Micheyl C (2011) Temporal coherence and attention in auditory

609        scene analysis. Trends Neurosci 34:114-123.

610    Shinn-Cunningham BG (2008) Object-based auditory and visual attention. Trends Cogn

611        Sci 12:182-186.

612    Steinschneider M, Liégeois-Chauvel C, Brugge JF (2011) Auditory evoked potentials and

613        their utility in the assessment of complex sound processing. In: The auditory cortex,

614        pp 535-559: Springer.

615    Steinschneider M, Nourski KV, Rhone AE, Kawasaki H, Oya H, Howard MA, 3rd

616        (2014) Differential activation of human core, non-core and auditory-related cortex

32

617   during speech categorization tasks as revealed by intracranial recordings. Front

618   Neurosci 8:240.

619  Sussman ES, Bregman AS, Wang WJ, Khan FJ (2005) Attentional modulation of

620   electrophysiological activity in auditory cortex for unattended sounds within

621   multistream auditory environments. Cogn Affect Behav Neurosci 5:93-110.

622  Sweet RA, Dorph-Petersen KA, Lewis DA (2005) Mapping auditory core, lateral belt,

623   and parabelt cortices in the human superior temporal gyrus. J Comp Neurol

624   491:270-289.

625  Yuille A, Kersten D (2006) Vision as Bayesian inference: analysis by synthesis? Trends

626   in cognitive sciences 10:301-308.

627  Zion Golumbic E, Cogan GB, Schroeder CE, Poeppel D (2013a) Visual input enhances

628   selective speech envelope tracking in auditory cortex at a "cocktail party". J

629   Neurosci 33:1417-1426.

630  Zion Golumbic EM, Poeppel D, Schroeder CE (2012) Temporal context in speech

631   processing and attentional stream selection: a behavioral and neural perspective.

632   Brain Lang 122:151-161.

633  Zion Golumbic EM, Ding N, Bickel S, Lakatos P, Schevon CA, McKhann GM,

634   Goodman RR, Emerson R, Mehta AD, Simon JZ, Poeppel D, Schroeder CE (2013b)

33

635        Mechanisms underlying selective neuronal tracking of attended speech at a "cocktail

636        party". Neuron 77:980-991.

637

34

638    **Legend:**

639    Figure 1: Illustrations of different decoding- and encoding-based neural

640    representations of the auditory scene and its constituents. (*Left)* Examples of

641    predicted MEG neural response using the Early-late model (red) and the

642    Summation model (magenta) superimposed on actual MEG response (black). The

643    proposed Early-late model prediction shows higher correlation with the actual

644    MEG neural response than Summation model. (*Right)* Example of speech

645    envelopes reconstructed (grey) from their late neural responses, for both the

646    foreground and the background, superimposed on actual speech envelopes of

647    foreground (blue) and background (cyan). The foreground reconstruction shows

648    higher correlation with the actual foreground envelope, compared to the

649    background reconstruction with the actual background envelope. All examples are

650    grand averages across subjects (3 seconds duration).

651

652    Figure 2: Early vs. late MEG neural responses to a continuous speech stimulus. A

653    sample stimulus envelope and multi-channel MEG recordings are shown in red

654    and black respectively. The two grey vertical lines indicate two arbitrary time

655    points at $t - \Delta t$ and $t$. The dashed and dotted boxes represent the early and late

656    MEG neural responses to stimulus at time point $t$ respectively. The reconstruction

35

657    of the stimulus envelope at time $t$ can be based on either early or late neural

658    responses, and the separate reconstructions can be compared against each other.

659

660

661    Figure 3: Stimulus envelope reconstruction accuracy using *early* neural responses.

662    **A.** Scatter plot of reconstruction accuracy of the foreground vs. individual

663    background envelopes. No significant difference was observed ($p = 0.21$), and

664    therefore no preferential representation of the foreground speech over the

665    individual background streams is revealed in early neural responses. **B.** Scatter

666    plot of reconstruction accuracy of the envelope of the entire acoustic scene vs. that

667    of the sum of the envelopes of all three individual speech streams. The acoustic

668    scene is reconstructed more accurately (visually, most of data points fall above the

669    diagonal) as a whole than as the sum of individual components in early neural

670    responses ($p < 2 \times 10^{-6}$). Reconstruction accuracy is measured by proportion of the

671    variance explained: the square of the Pearson correlation coefficient between the

672    actual and predicted envelopes.

673

674    Figure 4: Stimulus envelope reconstruction accuracy using *late* neural responses.

675    **A.** Scatter plot of accuracy between foreground vs. individual background

676    envelope reconstructions demonstrates that the foreground is represented with
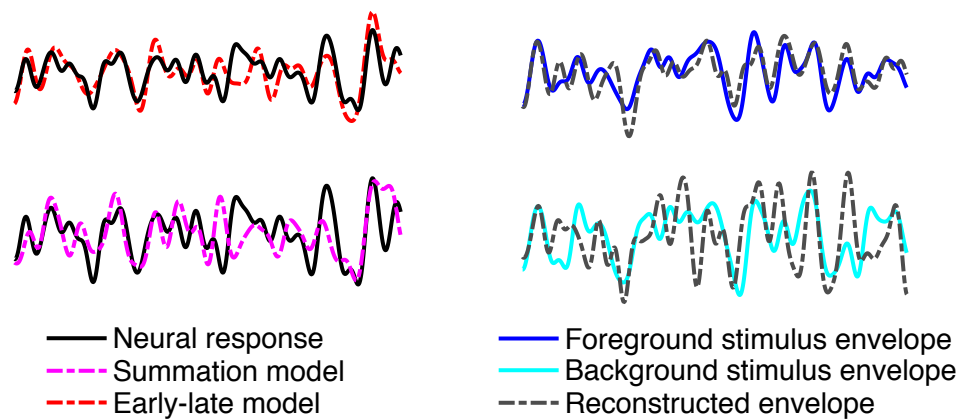
36

677    dramatically better fidelity (visually, most of data points fall above the diagonal)

678    than the background speech, in late neural responses ($p < 2 \times 10^{-6}$). **B.** Scatter plot

679    of the reconstruction accuracy of the envelope of the entire background vs. that of

680    the sum of the envelopes of the two individual background speech streams. The

681    background scene is reconstructed more accurately as a monolithic background

682    than as separated individual background streams in late neural responses ($p =$

683    0.012)

684

685    Figure 5: MEG response prediction accuracy. Scatter plot of the accuracy of

686    predicted MEG neural response for the proposed Early-late model vs. the standard

687    Summation model. The Early-late model predicts the MEG neural response

688    dramatically better (visually, most of data points fall above the diagonal) than the

689    Summation model ($p < 2 \times 10^{-6}$). The accuracy of predicted MEG neural

690    responses is measured by proportion of the variance explained: the square of the

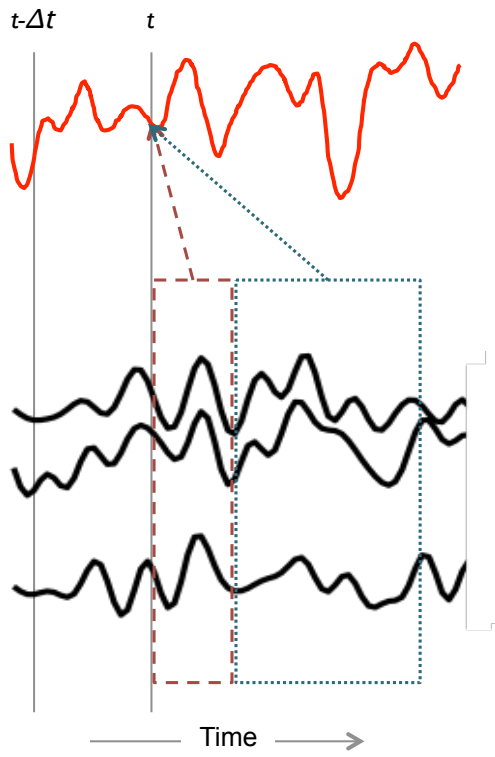691    Pearson correlation coefficient between the actual and predicted responses.

692

693

694

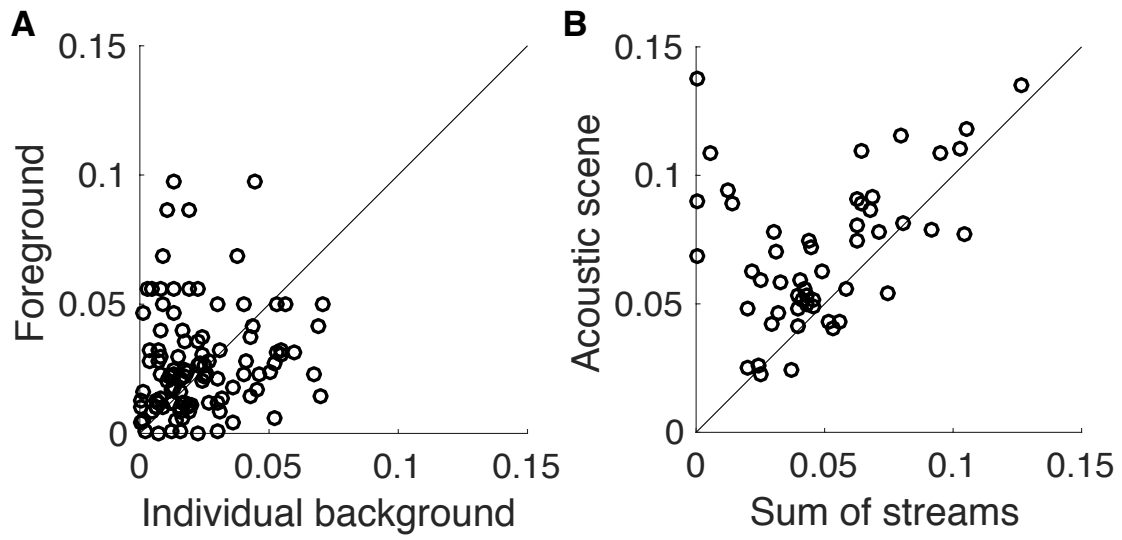695     Figure 1

38

696

697

698     Figure 2

39

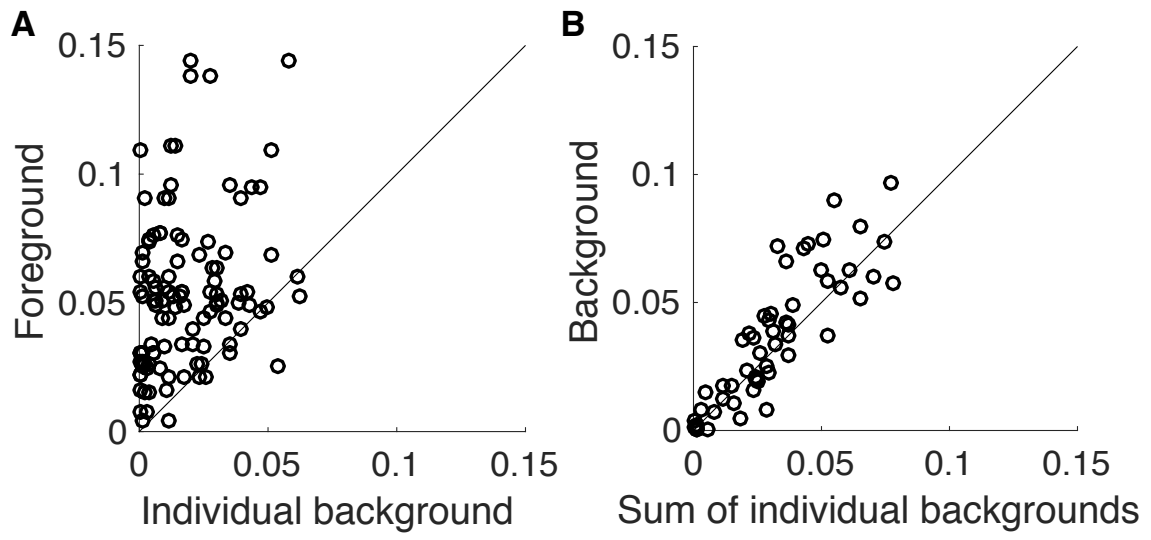# Stimulus Reconstruction Accuracy from **Early** Neural Responses
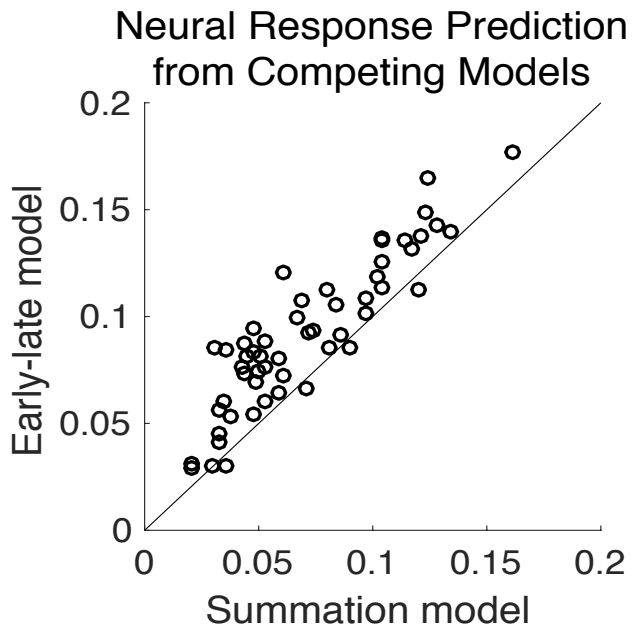


699

700

701    Figure 3

# Stimulus Reconstruction Accuracy from **Late** Neural Responses



702

703

704   Figure 4

Neural Response Prediction from Competing Models

705

706

707    Figure 5

42