

1 **Joint population coding and temporal coherence link an attended talker's voice and location**  
2 **features in naturalistic multi-talker scenes**

3 **Authors:** Kiki van der Heijden<sup>1,2,3</sup>, Prachi Patel<sup>2,4</sup>, Stephan Bickel<sup>5,6</sup>, Jose L. Herrero<sup>5,6</sup>, Ashesh D. Mehta<sup>5,6</sup>,  
4 Nima Mesgarani<sup>2,4,\*</sup>.

5 **Affiliations:**

- 6 1. Donders Institute for Brain, Cognition and Behavior, Radboud University, Nijmegen,  
7 Netherlands.
- 8 2. Mortimer B. Zuckerman Mind Brain Behavior Institute, Columbia University, New York, United  
9 States.
- 10 3. Maastricht Center for Systems Biology (MaCSBio), Maastricht University, Maastricht,  
11 Netherlands.
- 12 4. Department of Electrical Engineering, Columbia University, New York, United States.
- 13 5. Hofstra Northwell School of Medicine, New York City, United States.
- 14 6. The Feinstein Institute for Medical Research, New York City, United States.

15 **HIGHLIGHTS**

- 16 • Cortical responses to an single talker exhibit a distributed gradient, ranging from sites that are  
17 sensitive to both a talker's voice and location (dual-feature sensitive sites) to sites that are  
18 sensitive to either voice or location (single-feature sensitive sites).
- 19 • Population response patterns of dual-feature sensitive sites encode voice and location features  
20 of the attended talker in multi-talker scenes jointly and with equal precision.
- 21 • Despite their sensitivity to a single feature at the level of individual cortical sites, population  
22 response patterns of single-feature sensitive sites also encode location and voice features of a  
23 talker jointly, but with higher precision for the feature they are primarily sensitive to.
- 24 • Neural sites which selectively track an attended speech stream concurrently encode the  
25 attended talker's voice and location features.
- 26 • Attention selectively enhances temporal coherence between voice and location selective sites  
27 over time.
- 28 • Joint population coding as well as temporal coherence mechanisms underlie distributed multi-  
29 dimensional auditory object encoding in auditory cortex.

30 **ABSTRACT (240 words)**

31 Listeners readily extract multi-dimensional auditory objects such as a 'localized talker' from complex  
32 acoustic scenes with multiple talkers. Yet, the neural mechanisms underlying simultaneous encoding and

33 linking of different sound features – for example, a talker’s voice and location – are poorly understood.  
34 We analyzed invasive intracranial recordings in neurosurgical patients attending to a localized talker in  
35 real-life cocktail party scenarios. We found that sensitivity to an individual talker’s voice and location  
36 features was distributed throughout auditory cortex and that neural sites exhibited a gradient from  
37 sensitivity to a single feature to joint sensitivity to both features. On a population level, cortical response  
38 patterns of both dual-feature sensitive sites but also single-feature sensitive sites revealed simultaneous  
39 encoding of an attended talker’s voice and location features. However, for single-feature sensitive sites,  
40 the representation of the primary feature was more precise. Further, sites which selective tracked an  
41 attended speech stream concurrently encoded an attended talker’s voice and location features,  
42 indicating that such sites combine selective tracking of an attended auditory object with encoding of  
43 the object’s features. Finally, we found that attending a localized talker selectively enhanced temporal  
44 coherence between single-feature voice sensitive sites and single-feature location sensitive sites,  
45 providing an additional mechanism for linking voice and location in multi-talker scenes. These results  
46 demonstrate that a talker’s voice and location features are linked during multi-dimensional object  
47 formation in naturalistic multi-talker scenes by joint population coding as well as by temporal coherence  
48 between neural sites.

## 49 **SIGNIFICANCE STATEMENT**

50 Listeners effortlessly extract auditory objects from complex acoustic scenes consisting of multiple sound  
51 sources in naturalistic, spatial sound scenes. Yet, how the brain links different sound features to form a  
52 multi-dimensional auditory object is poorly understood. We investigated how neural responses encode  
53 and integrate an attended talker’s voice and location features in spatial multi-talker sound scenes to  
54 elucidate which neural mechanisms underlie simultaneous encoding and linking of different auditory  
55 features. Our results show that joint population coding as well as temporal coherence mechanisms  
56 contribute to distributed multi-dimensional auditory object encoding. These findings shed new light on  
57 cortical functional specialization and multidimensional auditory object formation in complex, naturalistic  
58 listening scenes.

## 59 **INTRODUCTION**

60 In everyday life, listeners rapidly and effortlessly parse complex acoustic scenes with multiple sound  
61 sources into its individual constituents. This process of auditory scene analysis (ASA<sup>1</sup>) is based on the  
62 segregation and subsequent grouping of features of temporally overlapping sound sources, resulting in  
63 the formation of coherent auditory objects<sup>2</sup>. Sound features contributing to auditory object formation  
64 include voice features related to object identity (e.g., pitch or timbre) and location features (e.g.,

65 interaural time differences, location cues)<sup>3-5</sup>. However, the neural basis for multi-dimensional auditory  
66 object formation in complex, naturalistic listening scenes is poorly understood.

67 One unresolved question is how cortical representations of individual sound features are linked by the  
68 brain to form a multi-dimensional auditory object. If voice and location features are encoded  
69 independently in two separate, functionally specialized and hierarchical processing streams as posited  
70 by the prevailing dual-stream framework<sup>6,7</sup>, it is not clear how these features are subsequently integrated  
71 to form a multi-dimensional auditory object. In contrast, recent studies using an active task design  
72 indicate that sound feature encoding may be distributed across auditory cortex rather than taking place  
73 in dedicated, functionally specialized anatomical regions as posited by the dual-stream theory. For  
74 example, studies in cats<sup>8</sup> and humans<sup>9</sup> showed that spatial sensitivity in primary auditory cortex (PAC)  
75 sharpens during goal-directed sound localization, suggesting that regions that are not considered part  
76 of the location pathway (i.e. PAC) may be recruited flexibly for spatial processing based on behavioral  
77 goals. Additionally, while speech processing has been attributed mostly to posterior STG<sup>10,11</sup>, a recent  
78 study demonstrated that speech processing is instead distributed across auditory cortex<sup>12</sup>. Such findings  
79 indicate that sound (feature) encoding may be more distributed than posited by the hierarchical dual-  
80 stream framework.

81 Additionally, it is not understood what neural mechanisms integrate cortical representations of individual  
82 sound features (e.g. spatial and non-spatial features). One hypothesis is that neuronal populations are  
83 sensitive to specific combinations of features and thereby encode multiple dimensions of an auditory  
84 object. Prior studies confirmed that some cortical sites are sensitive to multiple sound features  
85 simultaneously (e.g. in ferrets<sup>13</sup>, for a review<sup>14</sup>). However, because most prior measurements were  
86 performed with single sound sources, it is not known whether these cortical sites maintain their multi-  
87 dimensional sensitivity when presented with complex acoustic scenes comprising multiple, interfering  
88 sound sources. An alternative hypothesis states that auditory streams (pertaining to auditory objects)  
89 are formed through temporal coherence, i.e., response synchronization between neural populations that  
90 are sensitive to specific sound features<sup>15</sup>. Neural measurements in animals<sup>16,17</sup> and humans<sup>18</sup>  
91 demonstrate that temporal coherence is a plausible mechanism for auditory feature binding and  
92 segregation. It remains to be evaluated whether temporal coherence also underlies linking of voice and  
93 location features in human auditory cortex in naturalistic listening scenes.

94 Finally, although it is well known that auditory attention modulates the neural representation of spatial  
95 and non-spatial features<sup>19,20</sup> as well as auditory object formation<sup>21,22</sup>, it is not known how attention  
96 modulates integrated encoding of spatial and non-spatial features in complex, naturalistic sound scenes.  
97 Moreover, it remains an open debate<sup>2</sup> whether auditory objects form pre-attentively<sup>23</sup> or whether  
98 attention is necessary for auditory object formation<sup>15</sup>.

99 Here, we investigated cortical multi-dimensional auditory object formation with stereotactic  
100 electroencephalography (sEEG) recordings in neurosurgical patients. We measured neural activity in  
101 response to real-world sound scenes consisting of a single localized talker or two spatially separated  
102 talkers. The unique spatiotemporal resolution of neurophysiological recordings enabled us to map  
103 feature encoding and multi-dimensional object formation across auditory cortex. We found that active  
104 listening to complex, naturalistic scenes gives rise to distributed but joint voice and location encoding  
105 in single- as well as in multi-talker scenes. Furthermore, our results revealed that response patterns of  
106 distinct neural populations jointly encoded an attended talker's voice and location features. Finally, we  
107 show that attending to a localized talker in multi-talker scenes selectively enhanced temporal coherence  
108 between voice and location sensitive sites. In sum, these data demonstrate that multiple neural  
109 mechanisms contribute to linking an attended talker's voice and location in multi-talker scenes. .

## 110 **RESULTS**

111 We analyzed neural measurements in seven neurosurgical patients recorded with intracranial depth  
112 electrodes (stereoelectroencephalography, sEEG; Methods). Participants listened to English speech  
113 utterances consisting of one or two spatialized talkers. In single-talker scenes, either a male or female  
114 talker was present at a location of  $-45^\circ$  or  $+45^\circ$ . In two-talker scenes, a male and female talker were  
115 simultaneously present, one at  $-45^\circ$  and the other at  $+45^\circ$  (Figure 1 A). Trials had an average duration of  
116 5 s and the location of the talkers changed at random after each trial. The total duration of each condition  
117 (i.e., single-talker speech and multi-talker speech) was 8 minutes. For the single-talker condition, speech  
118 was paused at random intervals between trials and the participant was asked to repeat the last sentence  
119 as well as the location of the talker. For the multi-talker condition, participants were instructed at the  
120 start of a block to attend to a specific talker (i.e. 'attend male' or 'attend female'). At random moments  
121 in between trials, participants were asked to report the location of the attended talker and the last  
122 sentence uttered by the attended talker. Participants successfully performed the behavioral task (see <sup>24</sup>  
123 for a detailed analysis of the behavioral results).

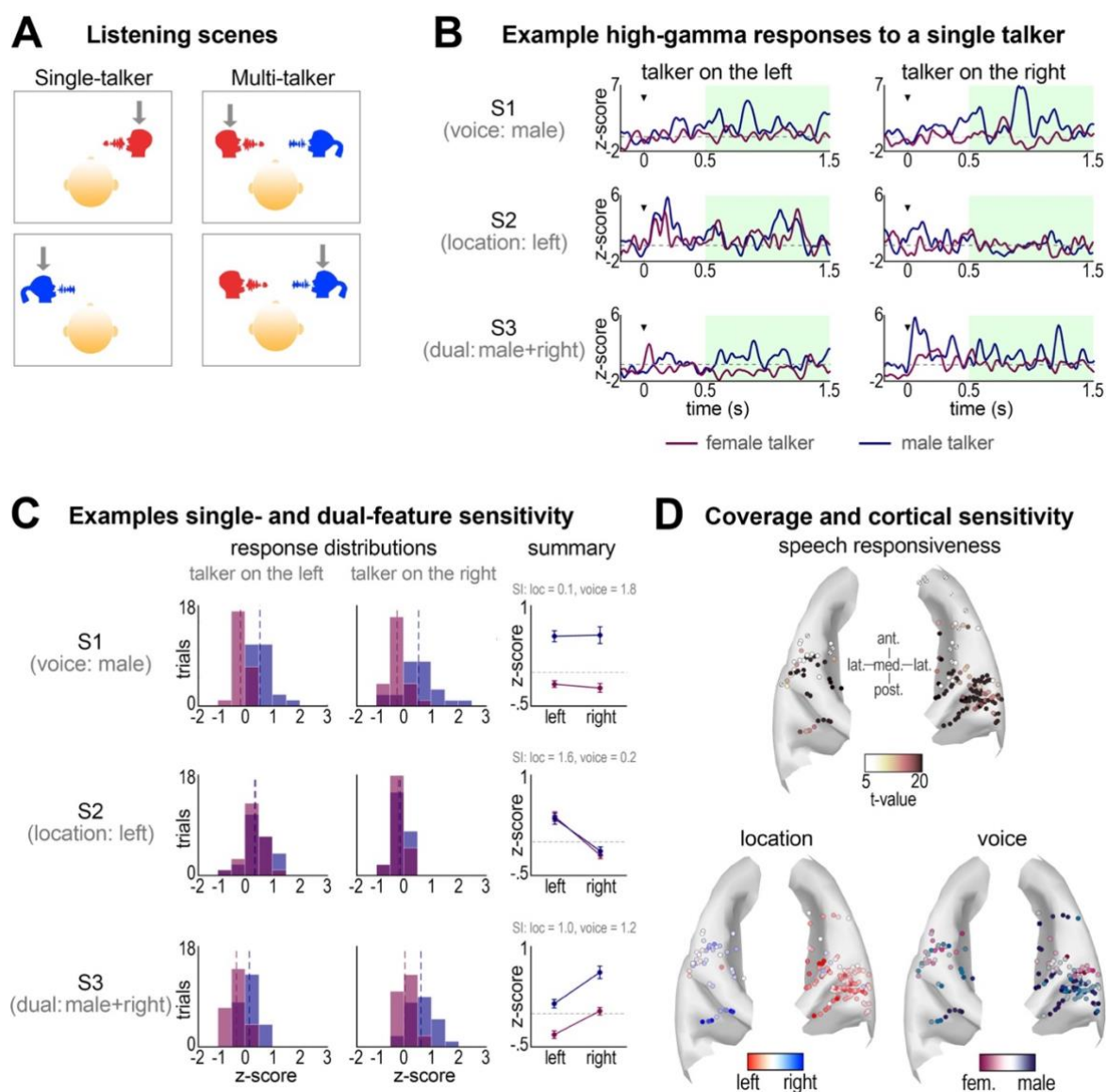
### 124 **Cortical sensitivity to a talker's voice and location features**

125 We observed significant neural population responses to speech in the high gamma envelope of 147  
126 cortical sites in auditory cortex (paired samples t-test of responses to speech versus silence,  $p < 0.05$ ,  
127 FDR corrected,  $q < 0.05$ ; Figure 1 D). These speech responsive sites were located in Heschl's gyrus (HG,  
128 6 left hemisphere, 32 right hemisphere), planum temporale (PT, 11 left hemisphere, 24 right hemisphere)  
129 and superior temporal gyrus (STG, 25 left hemisphere, 49 right hemisphere).

130 We characterized response properties for voice and location features by examining the responses to the  
131 single talker scenes for each cortical site. To assess to what extent a site exhibited sensitivity to voice, to

132 location, or to both, we contrasted the responses to one class of a feature (e.g., the male voice) to the  
133 responses to the other class of the feature (e.g., the female voice). For all sites, we extracted the mean  
134 response for each trial as the mean from 0.5 s post sound onset to 1.5 s post sound onset (that is,  
135 excluding the onset response). Figure 1 B shows example neural responses of three sites: One site  
136 sensitive to voice features (top panels), one site sensitive to location features (middle panels) and one  
137 site sensitive to both voice and location features (bottom panels). Figure 1 C shows the resulting  
138 response distributions for the sites in Figure 1 B. To test for sensitivity to voice features, we computed  
139 the effect size (Cohen's  $d^{p1}$ ) for the difference between the mean responses to all male and female trials,  
140 irrespective of the location of the talker (50 trials each). To test for sensitivity to location features, we  
141 computed Cohen's  $d$  for the difference between the mean responses to all trials in which the talker was  
142 at the right and all trials in which the talker was at the left, irrespective of the talker's voice (50 trials  
143 each). Figure 1 D depicts voice and location sensitivity (Cohen's  $d$ ) on the cortical surface. There was no  
144 overall relationship between sensitivity strength for a single talker's voice and location features ( $|$ Cohen's  
145  $d|$ ,  $r = 0.037$ ,  $p = 0.66$ ).

146 Statistical testing confirmed that 47 sites were significantly sensitive to voice features only (paired  
147 samples t-tests,  $p < 0.05$ , FDR corrected) and 12 sites were significantly sensitive to location features  
148 only (paired samples t-tests,  $p < 0.05$ , FDR corrected). In agreement with prior results (e.g. <sup>24,25</sup>), most  
149 sites which were sensitive to location, preferred locations in the contralateral hemifield. Further, 23 sites  
150 were sensitive to both location and voice features ( $p < 0.05$  for both t-tests). While multi-dimensional  
151 sensitivity has only been demonstrated for combinations of non-spatial features in humans, these results  
152 confirm prior work in animals<sup>26</sup> which showed that some neuronal populations in auditory cortex are  
153 sensitive for both spatial and non-spatial features<sup>13</sup>. In sum, cortical responses reveal a gradient from  
154 single-feature voice or location sensitive sites to dual-feature voice and location sensitive sites.



**Figure 1. Experiment design and single-talker cortical sensitivity for location and voice features.**

(A) Two examples of single-talker scenes (left panels) and two examples of spatial multi-talker scenes (right panels). Gray arrows indicate the attended talker. (B) Example neural responses from three sites: a single-feature voice sensitive site (S1, higher responses to male talker than to female talker irrespective of location), a single-feature location sensitive site (S2, higher responses to a talker on the left than a talker on the right, irrespective of the talker), and a dual-feature sensitive site (S3, higher responses to a male talker than to a female talker and higher responses to a talker on the right than a talker on the left). Black triangle indicates sound onset. Shaded green area depicts the time window for calculating voice and location sensitivity (i.e., 500 – 1,500 ms post sound onset). (C) Distribution of average trial responses to the male and female talker for three example sites (same as in B). Dashed line indicates the median of each distribution. Panels on the right depict the mean and standard error of the mean for each distribution. The sensitivity index (SI) is the effect size of the difference in response to two locations (SI loc) or the difference in response to two talkers (SI voice). (D) Top panel: Speech responsiveness of all electrodes in AC. Color saturation reflects the  $t$ -value for the contrast speech versus silence (see Methods). Electrodes that did not exhibit a significant response to speech are indicated by a slanted black line. Lower panels: Sensitivity for a single talker's location (left panel)

and voice features (right panel) plotted on the cortical surface for all speech responsive sites. Color indicates Cohen's  $d$  (range [-1,1]).

## 155 **Spectrotemporal tuning properties explain sensitivity to a talker's voice and location features**

156 Prior work showed that spectrotemporal tuning properties explain preferential responses to a talker's  
157 voice<sup>21,27</sup>. We examined whether we observe a similar relationship between spectrotemporal tuning and  
158 sensitivity to a talker's voice features in the present dataset and, additionally, we examine to what extent  
159 spectrotemporal tuning properties can also explain sensitivity to a talker's location features. We  
160 characterize the spectrotemporal tuning properties of each speech responsive site by estimating a  
161 spectrotemporal receptive field (STRF) from the responses to single-talker stimuli. We estimated STRFs  
162 using a five-fold cross-validation procedure, leaving out 20 trials and fitting the STRF on the remaining  
163 80 trials. We used the left-out 20 trials to estimate the goodness of fit, calculating the correlation  
164 between these left-out neural responses and neural responses predicted by the fitted STRFs (Methods).  
165 Next, we examined to what extent STRF shape explained sensitivity to talker's voice and location features  
166 for all cortical sites with a well-fitted STRF (correlation  $r > 0.2$ ,  $n = 93$ ).

167 To analyze the relationship between STRF shape and sensitivity to a talker's voice features, we divided  
168 the group of sites sensitive to a talker's voice ( $n = 47$ , Figure 1) into sites responding maximally to the  
169 male talker and sites responding maximally to the female talker. In line with prior work<sup>21</sup>, the average  
170 STRF of sites responding preferentially to the female talker exhibited tuning properties corresponding  
171 to the spectral profile of the female talker, while the average STRF across sites responding preferentially  
172 to the 'male' talker exhibited tuning properties corresponding to the spectral profile of the male talker.  
173 That is, Figure 2 A shows that the average STRF of 'male'-preferring sites exhibited an excitatory region  
174 at low frequencies between 50 Hz and 100 Hz, overlapping with F0 of the male talker (65 Hz). In contrast,  
175 the average STRF of 'female'-preferring sites exhibited an excitatory region between 160 Hz and 200 Hz,  
176 overlapping with F0 of the female talker (175 Hz).

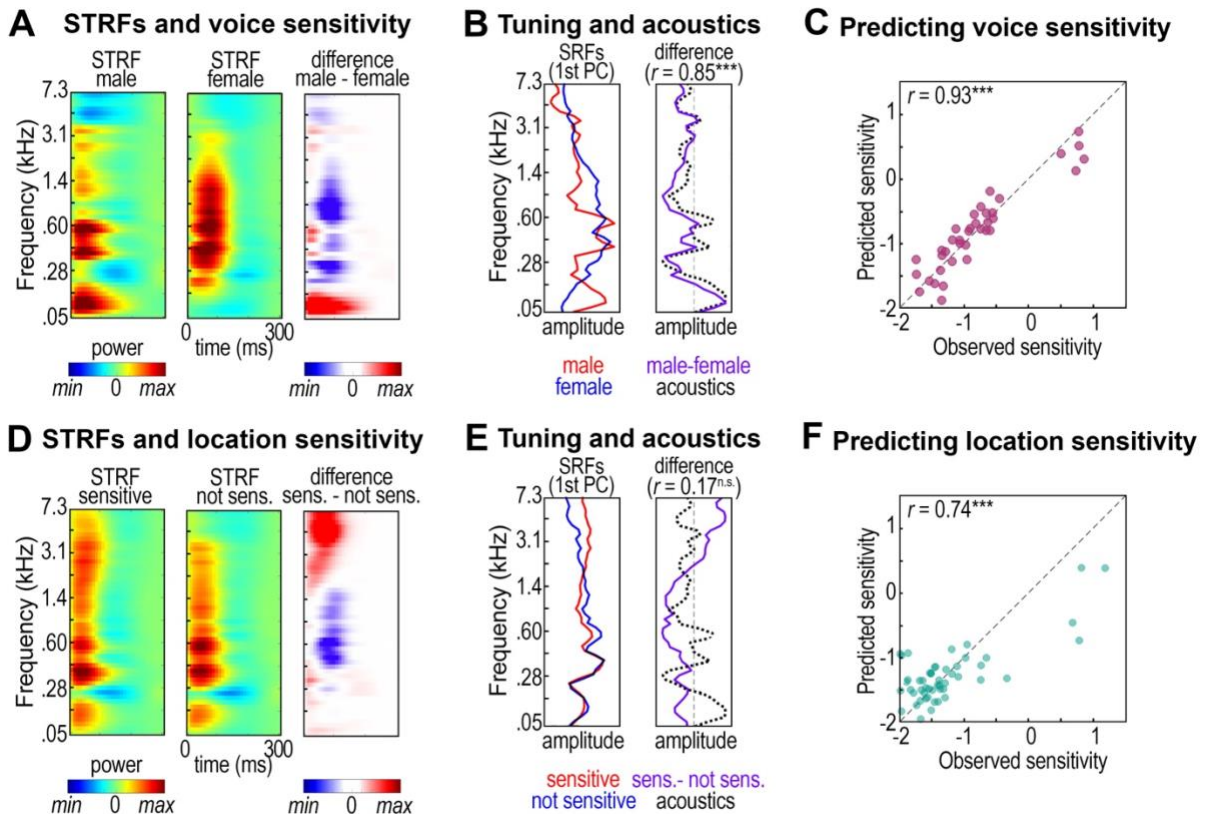
177 To quantify this relationship between spectral tuning properties and sensitivity to a talker's voice  
178 features, we extracted the spectral receptive fields (SRFs) of sites responding maximally to the male talker  
179 and sites responding maximally to the female talker. The SRF corresponds to the first component of a  
180 principal component analysis (PCA) of the STRF along the spectral dimension<sup>21</sup>. The difference in the  
181 SRFs of these two groups (i.e., responding preferentially to the female or male voice) was strongly  
182 correlated to the difference between the spectral profile of the male and female talker, indicating that  
183 the preference of voice sensitive sites for the male or female talker was driven by the correspondence  
184 between the spectral response profile of the site and the acoustic profile of the talker ( $r = 0.85$ ,  $p = 7.1E-$   
185  $15$ ; Methods; Fig. 2 B). Further, mapping the SRFs to sensitivity for a talker's voice using ridge regression

186 (Methods), showed that SRFs predicted sensitivity for a talker's voice well. That is, there was a high  
187 correlation between predicted sensitivity and observed sensitivity ( $r = 0.932$ ,  $p = 5.9E=17$ ; Fig. 2 C). These  
188 findings confirm that sensitivity to a talker's voice is driven by the spectral tuning properties of cortical  
189 sites<sup>21</sup>.

190 We then repeated the STRF analysis to assess the relationship between STRF shape and sensitivity to a  
191 talker's location features. First, we computed the average STRF across sites that are sensitive to a talker's  
192 location ( $n = 12$ , Figure 1) and across sites that were not sensitive to a talker's location (and not sensitive  
193 to a talker's voice either,  $n = 65$ ). Figure 2 D shows that sites which were sensitive to a talker's location  
194 had an excitatory STRF region for frequencies above 1.5 kHz which was reduced in sites which were not  
195 sensitive to a talker's location. In contrast, sites which were sensitive to a talker's location responded  
196 more weakly to frequencies between 0.5 – 1.5 kHz. No difference in STRF properties was observed for  
197 frequencies below 0.5 kHz.

198 As expected, extracting and comparing the SRFs of the two groups (i.e., location-sensitive versus not  
199 location-sensitive) showed that the difference in SRFs was not correlated to the difference between the  
200 spectral profile of the male and female talker ( $r = 0.17$ ,  $p = 0.23$ ; Methods; Fig. 2 D). However, mapping  
201 SRFs to sensitivity to a talker's location features (Methods), demonstrated that SRFs predicted such  
202 location-sensitivity well: There was a high correlation between predicted and observed sensitivity ( $r =$   
203  $0.932$ ,  $p = 5.9E=17$ ; Fig. 2 C). These findings indicate that sites which are sensitive to a talker's location  
204 respond more strongly to frequencies with robust interaural level difference (ILD) cues for sound  
205 localization<sup>28</sup>, while responding less strongly to frequencies in which binaural disparity cues such as ILDs  
206 and interaural time differences (ITDs) are less reliable<sup>29</sup>. Further, these results indicate that sensitivity to  
207 a talker's location is not related to tuning to low frequencies (i.e.,  $< 0.5$  kHz). Taken together, these  
208 findings confirm previous work demonstrating that spectrotemporal tuning explains tuning to a talker's  
209 voice<sup>21,27</sup> and extend this by showing that spectrotemporal tuning also explains tuning to a talker's  
210 location.





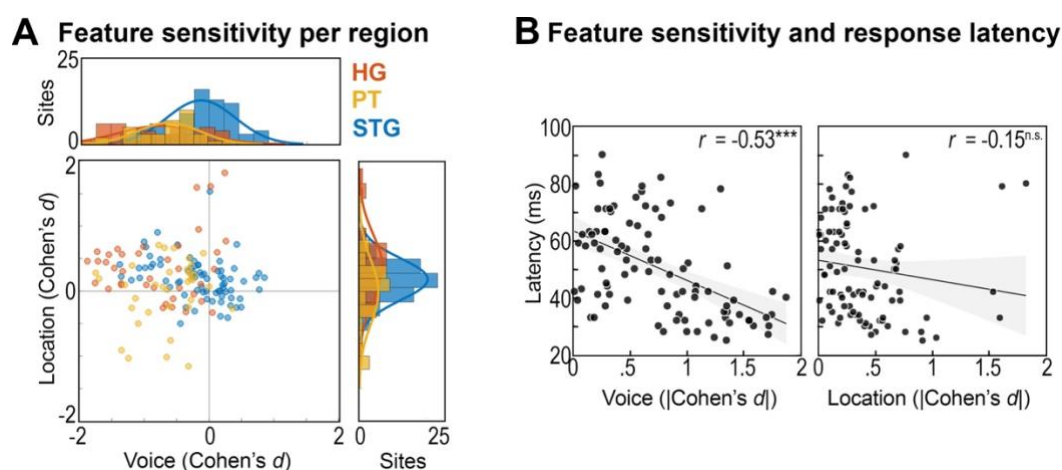
**Figure 2. Spectrotemporal tuning characteristics explain sensitivity to a talker's voice and to a talker's location.** (A) Spectrotemporal tuning properties related to voice sensitivity. From left to right: Average STRF for sites responding maximally to the male talker, average STRF for sites responding maximally to the female talker and the difference (STRF male – STRF female). (B) Comparing spectral tuning properties to the acoustics of the male and female talker. Left panel: Average spectral receptive field of sites responding maximally to a female talker (blue). Right panel: The correlation between the difference SRF (SRF male – SRF female) and the difference in the acoustics of the male and female talker. (C) Predicting voice sensitivity from the difference SRF (Cohen's  $d$ ). (D) Spectrotemporal tuning properties related to location sensitivity. From left to right: Average STRF for location sensitive sites, average STRF for sites that were not sensitive to location and the difference (STRF sensitive – not sensitive). (E) Comparing spectral tuning properties to location sensitivity. Left panel: Average spectral receptive field of sites sensitive to location features (red) and for sites not sensitive to location features (blue). Right panel: No correlation between the difference SRF (SRF not location sensitive – SRF location sensitive) and the difference in the acoustics of the male and female talker. (F) Predicting location sensitivity (Cohen's  $d$ ) from the difference SRF.

## 211 Sensitivity to a talker's voice and location across the cortical hierarchy

212 To investigate to what extent sensitivity to a talker's voice and location can be related to cortical  
 213 processing stages, we investigated how sensitivity to a talker's features was distributed across auditory  
 214 cortex. While several studies linked delineated anatomical regions to hierarchical processing stages (for  
 215 example, HG is considered primary auditory cortex and PT and STG higher-order auditory regions<sup>30</sup>),  
 216 other work investigating neural response latencies and response properties showed that a single  
 217 anatomical region may contain different auditory processing stages (e.g. <sup>12,31</sup>). That is, as response

218 latency roughly corresponds to the number of synapses away from the periphery it is considered as an  
219 indication of the processing stage of a neural site. Here, we therefore assessed the distribution of feature  
220 sensitivity both within cortical auditory regions and as a function of response latency. We calculated  
221 response latency as the peak along the temporal dimension of the STRF (for sites with a well-fitted STRF,  
222  $r > 0.2$ ,  $n = 93$ ; Methods).

223 Figure 3 A shows the regional distributions of Cohen's  $d$  for a talker's voice. Comparing the distributions  
224 showed that sensitivity to a talker's voice was stronger in HG than in STG ( $|\text{Cohen's } d|$ , Kruskal-Wallis H  
225 test,  $\chi^2(2) = 14.6$ ,  $p = 0.0007$ , Figure 3 A). Further, there was a negative correlation between sensitivity  
226 to a talker's voice and response latency ( $r = -0.526$ ,  $p = 1.3\text{E-}7$ ; Figure 3 B). These findings confirm prior  
227 reports of a decrease in sensitivity to a talker's voice along the cortical auditory processing hierarchy<sup>21</sup>.  
228 In contrast, although we observed a trend towards regional differences in the distribution of Cohen's  $d$   
229 for a talker's location ( $|\text{Cohen's } d|$ , Kruskal-Wallis H test,  $\chi^2(2) = 5.45$ ,  $p = 0.07$ ; Figure 3 A), this trend  
230 failed to reach significance. Moreover, we observed no correlation between sensitivity to a talker's  
231 location and response latency ( $r = -0.146$ ,  $p = 0.16$ ; Figure 3 B). While the lack of regional differences  
232 may be a consequence of the relatively low anatomical sampling density, together the regional and  
233 response latency results indicate that sensitivity is consistent across low- and high-level processing  
234 stages during active listening. These findings confirm recent work<sup>9,24</sup>, but contrast the predictions of the  
235 dual-stream framework which posits that PT is functionally specialized for spatial processing<sup>6,7,32</sup>.



**Figure 3. Sensitivity to a single talker's voice and location across the cortical hierarchy** (A) Scatterplot of voice sensitivity (x-axis) and location sensitivity (y-axis). Each symbol represents an individual site. Bar graphs depict corresponding marginal distributions for voice sensitivity (left) and location sensitivity (right). (B) Correlation between single-talker response latency and feature sensitivity (left panel: voice; right panel: location). Each circle depicts a site. Solid lines depict the correlation; shaded areas depict the 95% confidence interval. Asterisks indicate significance:  $*** = p < 0.001$ .

236

## 237 **Attentional modulation of neural responses to a talker's voice and location in multi-talker scenes**

238 We showed that cortical sites exhibit varying degrees of sensitivity for a single talker's location and voice.  
239 Motivated by prior findings of attentional modulation of neural responses to a talker's voice and  
240 location<sup>21,24</sup>, we examined to what degree we observed such local attentional modulations in multi-talker  
241 scenes in our data. Further, we characterized how attentional modulation by a talker's voice relates to  
242 attention modulation by a talker's location. Specifically, we quantified to what extent attending to a  
243 talker's voice and location in multi-talker scenes modulated the response gain of individual cortical sites  
244 similar to our quantification of sensitivity to a single talker's voice and location features. Specifically, we  
245 calculated the effect size Cohen's  $d$  for the difference in mean response to the trials for each attentional  
246 condition. As before, we calculated the mean response for each trial from 0.5 s post sound onset to 1.5  
247 s post sound onset (excluding the onset response).

248 In agreement with prior studies<sup>21,24</sup>, attending a localized talker evoked weak response gain modulations  
249 across speech responsive sites both by the attended talker's voice and by the attended talker's location.  
250 Figure 4 A shows that attentional modulation of response gain was smaller than modulation by a single  
251 talker's voice or location, that is, single-talker sensitivity (paired samples t-test of |Cohen's  $d$ |; voice:  
252  $t(146) = 9.65$ ,  $p = 2.38E-17$ ; location:  $t(146) = 6.68$ ,  $p = 4.64E-10$ ). In agreement with this, statistical  
253 testing did not identify neural sites of which the response gain was modulated significantly by attention  
254 to the talker's voice (paired samples t-tests,  $p > 0.05$ ), the talker's location (paired samples t-tests,  $p >$   
255  $0.05$ ), or jointly ( $p > 0.05$ ). Further, Figure 4 B shows that only few sites were jointly modulated by an  
256 attended talker's voice and location in multi-talker scenes. Specifically, only few electrodes were close  
257 to the diagonal and exhibited attentional response gain modulation for both voice and location  
258 ( $|Cohen's\ d| > .1$ ). Crucially, as single-source sensitivity to a specific sound feature (e.g. pitch, location) is  
259 generally considered an indication of functionally specialized processing<sup>2,32</sup>, the lack of corresponding,  
260 dedicated attentional modulations of response gains raises the question what the role is of these sites  
261 in the encoding of an attended sound source in scenes with multiple sound sources.

## 262 **Decoding a localized talker from population activity patterns in multi-talker scenes**

263 To elucidate the relationship between local encoding properties and population encoding properties,  
264 we examined whether a localized talker can be decoded from population response patterns. Specifically,  
265 we used a linear decoding approach to assess to what extent a localized talker can be decoded from  
266 population responses in single talker scenes and to what extent an attended localized talker can be  
267 decoded from population responses in multi-talker scenes. To decode a localized talker in single-talker  
268 scenes, we trained a four-class regularized least-squares (RLS<sup>22,33</sup>) classifier on the response patterns in  
269 single-talker scenes using a leave-two-trials-out cross-validation procedure (corresponding to 25 folds).

270 To decode an attended localized talker, we trained an identical four-class RLS classifier on the response  
271 patterns in multi-talker scenes using a similar cross-validation procedure. We assessed decoding  
272 accuracy by predicting the talker's voice and location from the response patterns of the left-out trials of  
273 each fold (Methods).

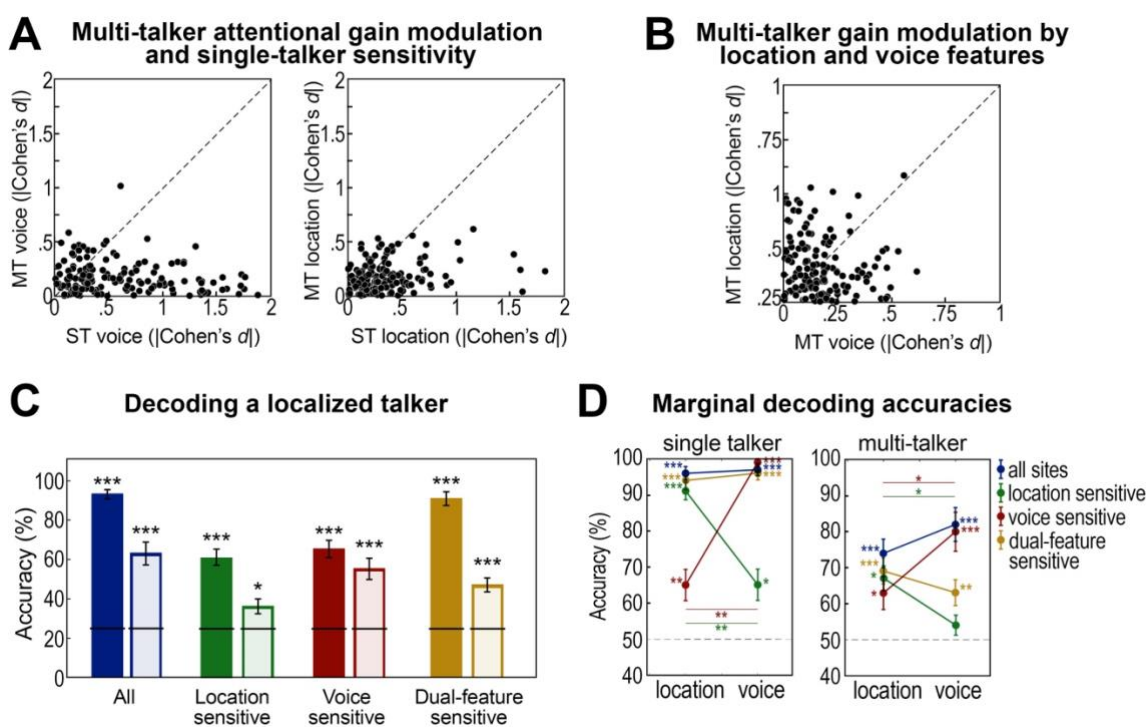
274 As expected, Figure 4 C shows that a single localized talker could be accurately decoded from the entire  
275 population of speech responsive sites ( $n = 147$ ; average accuracy [standard error of the mean; SEM] =  
276 93.0 % [2.29],  $p = 0$ , FDR corrected). Similarly, the attended localized talker was decoded accurately from  
277 the entire population of speech responsive sites (mean accuracy [SEM] = 63.0 % [5.80],  $p = 0$ ). Marginal  
278 decoding accuracies for the talker's voice and location show that both features were decoded with equal  
279 precision in single-talker scenes (Figure 4 D, mean marginal accuracy: voice [SEM] = 97.0 % [1.66],  $p =$   
280 0; location [SEM] = 96.0 % [1.87],  $p = 0$ ; paired samples  $t$ -test,  $t(24) = 0.37$ ,  $p = 0.72$ ) as well as in multi-  
281 talker scenes (Figure 3D, mean accuracy voice [SEM] = 82.0 % [4.68],  $p = 0$ ; average accuracy location  
282 [SEM] = 74.0 % [3.95],  $p = 0$ ;  $t(24) = 1.69$ ,  $p = 0.14$ ). These findings show that although attentional  
283 modulation of local response gain by the attended talker's voice and location in multi-talker scenes was  
284 weak (Figure 4 A), response patterns across the entire population of speech responsive sites the attended  
285 localized talker with high fidelity.

286 Next, we examined how sites which exhibit single-feature sensitivity for a talker's voice or location  
287 features in their local responses ( $n = 47$  and  $n = 12$ , Figure 1) encode a localized talker in population  
288 response patterns. We therefore trained the RLS classifier on the population responses of these neural  
289 sites with the same procedure described above. Note that although the latter population is relatively  
290 small, we chose to use this stringent selection to ensure that the population did not incorporate sites  
291 that were also to some extent sensitive to a talker's voice. Figure 4 C shows that the classifier successfully  
292 decoded a localized talker in single-talker scenes from voice sensitive sites (mean accuracy [SEM] = 65.0  
293 % [4.33],  $p = 0$ ) as well as from location sensitive sites (mean accuracy 61.0 % [4.10],  $p = 0$ ). Further, the  
294 marginal accuracies in Figure 4 D show that the classifier decoded the talker's voice more accurately  
295 from population responses of voice sensitive electrodes than the talker's location (mean marginal  
296 accuracy: voice [SEM] = 99.0 % [1.00],  $p = 0$ ; location [SEM] = 65.0 % [4.33],  $p = 0.0098$ ; paired samples  
297  $t$ -test,  $t(24) = 7.49$ ,  $p = 3.93E-7$ ). Conversely, decoding accuracy was higher for the talker's location than  
298 for the talker's voice when the classifier operated on population responses of location sensitive sites  
299 (voice [SEM] = 65.0 % [4.33],  $p = 0.02$ ; location [SEM] = 91.0 % [2.45],  $p = 0.0024$ ;  $t(24) = 5.32$ ,  $p = 3.74E-$   
300 5). These findings show that single-feature sensitive sites nevertheless encode coarse information about  
301 other feature dimensions of an individual talker in their population responses.

302 Furthermore, in multi-talker scenes, the classifier decoded the attended localized talker above chance  
303 level from population responses of voice sensitive sites (Figure 4 C, mean accuracy [SEM] = 55.0 % [5.59],

304  $p = 0$ ) as well as from population responses of location sensitive sites (Figure 4 C, mean accuracy [SEM]  
305 = 36.0 % [3.84],  $p = 0.026$ ). However, while both the attended talker's location and voice were decoded  
306 above chance level from population responses of voice sensitive sites (mean marginal accuracy: voice  
307 [SEM] = 80.0 % [5.40],  $p = 0$  ; location [SEM] = 63.0 % [4.59],  $p = 0.023$ ), only the attended talker's  
308 location was decoded accurately from population responses of location sensitive sites. Specifically, the  
309 decoding accuracy for the attended talker's voice just failed to reach statistical significance, which may  
310 be a consequence of the small number of sites in this group (mean marginal accuracy: voice [SEM] =  
311 54.0 % [2.77],  $p = 0.079$ ; location [SEM] = 67.0 % [3.45],  $p = 0.021$ ). For both voice and location sensitive  
312 sites, the preferred feature was decoded significantly better than the other feature (voice:  $t(24) = 2.72$ ,  $p$   
313 = 0.024; location:  $t(24) = 2.98$ ,  $p = 0.024$ ). These findings indicate that populations which exhibit local  
314 properties of functional specialization in response to single-source sound scenes may nonetheless  
315 encode (coarse) information about multiple dimensions of the auditory object. Future work including  
316 more fine-grained sampling of multiple feature dimensions (e.g., more talkers and more voices) is  
317 required to establish the resolution with which population responses of single-feature sensitive sites  
318 encode other feature dimensions.

319 Finally, we showed previously that some neural sites were sensitive both for a talker's voice and location  
320 (Figure 1,  $n = 23$ ). We examined to what extent population responses of these dual-feature sensitive  
321 sites also jointly encode a talker's voice and location in single- and in multi-talker conditions. That is,  
322 while prior work in animals showed that auditory cortex contains sites which are sensitive to spatial as  
323 well as non-spatial sound features<sup>13</sup>, work in humans focused only on multi-dimensional sensitivity for  
324 non-spatial features (e.g. <sup>26</sup>). Moreover, as all prior studies were conducted with single-source scenes, it  
325 is not clear to what extent multi-dimensional sensitivity is maintained in multi-talker sites. Here, Figure  
326 4 C shows that the localized talker was decoded accurately from population responses of dual-feature  
327 sensitive sites in single-talker scenes (mean accuracy [SEM] = 91.0 % [3.50],  $p = 0$ ). The attended localized  
328 talker was also decoded accurately from population responses in multi-talker scenes (mean accuracy  
329 [SEM] = 47.0 % [3.63],  $p = 0$ ). Importantly, Figure 4 D shows that dual-feature sensitive encoded the  
330 talker's voice and location with equal precision, both in single-talker (mean marginal accuracy: voice  
331 [SEM] = 96.0 % [1.87],  $p = 0$ ; location [SEM] = 94.0 % [2.61],  $p = 0$ ; paired samples t-test:  $t(24) = 0.81$ ,  $p$   
332 = 0.57) and in multi-talker scenes (mean marginal accuracy: voice [SEM] = 63.0 % [3.57],  $p = 0.006$ ;  
333 location [SEM] = 69.0 % [3.62],  $p = 0$ ; paired samples t-test:  $t(24) = 1.1$ ,  $p = 0.28$ ). In sum, we show that  
334 population responses of dual-feature sensitive sites encode both spatial and non-spatial features of an  
335 attended talker in multi-talker scenes. This suggests that such dual-feature sensitive sites contribute to  
336 the encoding of multiple dimensions of an auditory object.



**Figure 4. Attentional gain modulation and the representation of a localized talker in single- and multi-talker scenes.** (A) Relationship between attentional gain modulation and single-talker sensitivity. Left panel: Gain modulation by attending to a talker’s voice in multi-talker scenes (y-axis) versus single-talker voice sensitivity (x-axis). Right panel: Gain modulation by attending a talker’s location in multi-talker scenes (y-axis) versus single-talker location sensitivity (x-axis). (B) Scatterplot of gain modulation by an attended talker’s voice (x-axis, |Cohen’s  $d$ ) and by an attended talker’s location (y-axis, |Cohen’s  $d$ ). (C) Decoding a localized talker from response patterns in single-talker scenes (filled bars) and decoding an attended localized talker in multi-talker scenes (open bars). Horizontal lines depict chance level. Asterisks indicate significance: \*\*\* =  $p < 0.001$ . (D) Marginal decoding accuracies for a talker’s voice and location. Dashed line depicts chance level, asterisks indicate significance: \* =  $p < 0.05$ ; \*\* =  $p < 0.01$ ; \*\*\* =  $p < 0.001$ .

337 **Sites selectively tracking an attended speech stream simultaneously encode an attended talker’s**  
 338 **voice and location features**

339 In the preceding sections, we examined to what extent single-feature sensitive and dual-feature sensitive  
 340 sites encode an attended talker’s voice and location features in multi-talker scenes. However, prior work  
 341 showed that auditory cortex also contains neural sites which are not strongly sensitive to a single talker’s  
 342 features, but which nonetheless play a crucial role in auditory object formation by selectively tracking  
 343 the attended speech stream in multi-talker listening scenes<sup>21,22,24,34</sup>. As the relationship between such  
 344 speech stream tracking and encoding of the attended talker’s features is not known, we analyzed the  
 345 measured neural responses to multi-talker scenes to evaluate to what degree such neural sites which  
 346 selectively track sites additionally encode the attended talker’s voice and location features.

347 For each site, we first quantified selective tracking of the attended speech stream by calculating to what  
348 extent a site's responses in spatial multi-talker scenes were modulated by attention to reflect the  
349 response to the attended localized talker in single talker scenes. We define the tracking index (TI) for  
350 each site similar to the definition in <sup>21</sup>:

$$\begin{aligned} 351 \quad \text{Tracking Index} &= \text{corr}(att_M, single_M) - \text{corr}(att_M, single_F) + \text{corr}(att_F, single_F) \\ 352 \quad &\quad - \text{corr}(att_F, single_M) \end{aligned}$$

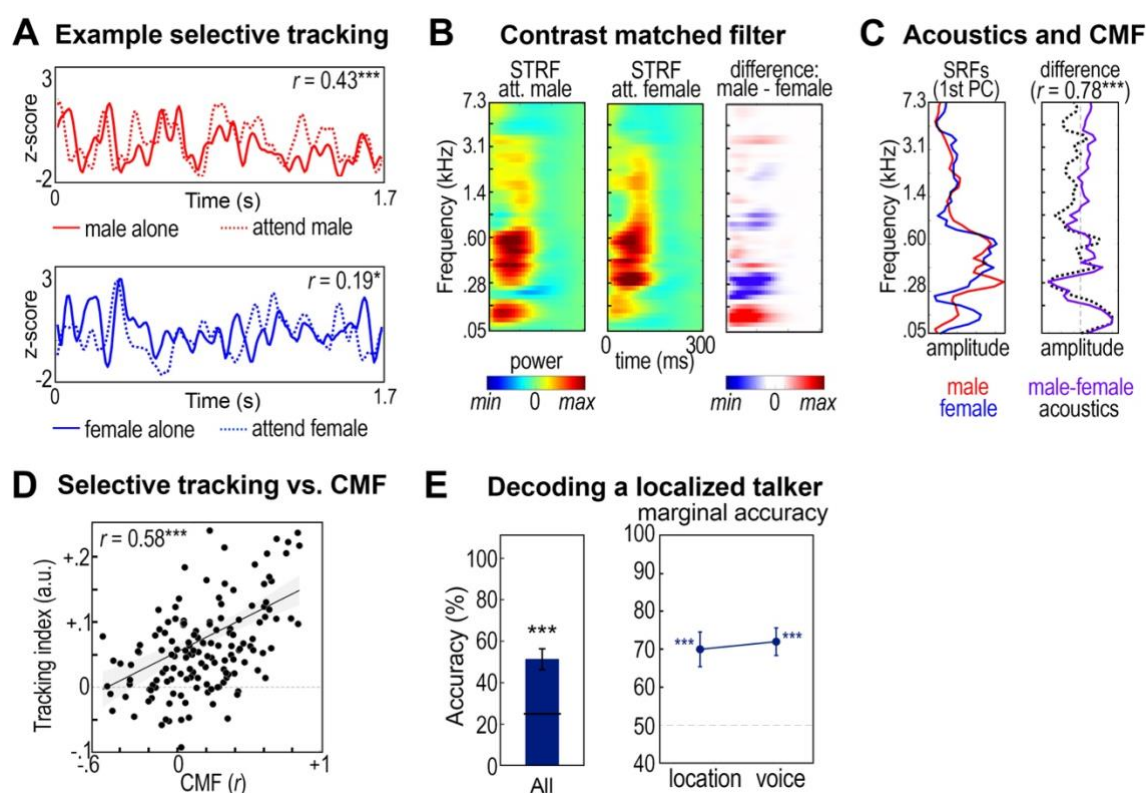
353 Here, *M* refers to the male talker and *F* to the female talker. Further,  $\text{corr}(att, single)$  corresponds to the  
354 correlation between the single-talker response and the multi-talker response for the same trial calculated  
355 over the entire duration over the trial (see example in Figure 5 A).

356 Next, to gain more insight into the encoding properties of sites which selectively track an attended  
357 speech stream, we examined to what extent TI is explained by attentional modulation of STRF  
358 properties<sup>35,36</sup>. We quantified such attentional modulation of STRF properties by estimating for each  
359 cortical site two STRFs from the responses to multi-talker scenes: one for the 'attend male talker'  
360 condition and one for the 'attend female talker' condition. To relate attention-induced spectrotemporal  
361 plasticity to the encoding of a talker's spectral characteristics, we extracted the spectral receptive field  
362 (SRF) from each STRF as the first principal component of a PCA (only for cortical sites with a robust STRF  
363 as estimated from single-talker responses,  $r > 0.2$ ,  $n = 93$ ) and compared these to the spectral profile of  
364 the male and female talker. Specifically, we computed the difference in SRF for the two attention  
365 conditions (i.e., attend female – attend male) and correlated this difference SRF to the acoustic difference  
366 spectrum between the male and female talker (see examples in Figure 5 B, C). If a site's spectral tuning  
367 properties are modulated by attention towards the attended talker's spectral profile, we expect a high  
368 correlation between the difference SRF and the acoustic difference spectrum (Fig. 5 C). That is, we expect  
369 these sites to resemble a contrast matched filter which facilitates figure-ground segregation by  
370 enhancing the attended target (e.g., the female talker) and filtering out the background (e.g., the male  
371 talker)<sup>37</sup>. Therefore, we quantified the strength of attentional modulation of STRF properties by  
372 calculating a contrast matched filter (CMF) index, which is the correlation between the attention-driven  
373 difference in the SRFs and the acoustic difference spectrum for the female and male talker (Figure 5 C).

374 As expected<sup>24</sup>, CMF explains TI well ( $r = 0.577$ ,  $p = 1.4E-9$ , Fig. 5 D). This indicates that sites whose STRF  
375 properties are strongly modulated by attention tend to be sites which selectively track the attended  
376 speech stream. In contrast, TI is not correlated to single talker encoding properties (voice sensitivity:  $r =$   
377  $0.11$ ,  $p = 0.24$ ; location sensitivity:  $r = 0.12$ ,  $p = 0.24$ ) or multi-talker attentional response gain  
378 modulation (attended talker's voice:  $r = 0.06$ ,  $p = 0.48$ ; attended talker's location:  $r = -0.13$ ,  $p = 0.24$ ).  
379 Thus, the encoding properties of sites which selectively track an attended speech stream are

380 characterized by attentional modulation of STRF properties rather than by single-talker sensitivity or  
 381 multi-talker attentional response gain modulation.

382 Further, to examine whether the population responses of sites which selectively track the attended  
 383 speech stream also encode the attended talker's voice and location, we trained the four-class classifier  
 384 on their population response patterns in multi-talker scenes (i.e, for sites with TI > 0.1,  $n = 33$ ). The  
 385 classifier accurately decoded the attended localized talker from these population response patterns  
 386 (average accuracy [SEM] = 51.0 % [5.10],  $p = 0$ ; Figure 5 E). Furthermore, the classifier decoded the  
 387 attended talker's voice and location with equal precision (Figure 5 E; marginal accuracies: voice [SEM] =  
 388 72.0 % [3.63],  $p = 0$ ; location [SEM] = 70.0 % [4.56],  $p = 0$ ; paired samples t-test,  $t(24) = 0.40$ ,  $p = 0.69$ ).  
 389 In sum, population responses of sites which selectively tracked the attended speech stream also encoded  
 390 the attended talker's voice and location. This finding indicates that the population responses of these  
 391 sites play a role in combining selective tracking of an attended auditory object (here, speech stream)  
 392 with encoding of the features of that object (here, the talker's voice and location).



**Figure 5. Selective speech tracking and encoding of the attended talker's voice and location.**

(A) High-gamma responses of an example site exhibiting selective tracking of the attended talker. Solid lines depict response in single-talker scene, dotted lines depict response in multi-talker scene. (B) STRFs of an example cortical site exhibiting contrast matched filtering. Left panel: STRF in the 'attend male' condition. Middle panel: STRF in the 'attend female' condition. Right panel: Difference (STRF attend male – STRF attend female). (C) Comparing spectral tuning properties in the two attention conditions to the acoustics of the male and female talker. Left panel: Spectral receptive fields



for the 'attend male' condition (red) and for the 'attend female' condition (blue). Right panel: The correlation between the difference SRF (SRF attend male – SRF attend female) and the difference in the acoustics of the male and female (D) Correlation between CMF (x-axis) and Tracking Index (y-axis). Circles represent cortical sites. (E) Left panel: Decoding an attended localized talker from population responses in multi-talker scenes. Horizontal line depicts chance level. Asterisks indicate significance: \*\*\* =  $p < 0.001$ . Right panel: Marginal decoding accuracies for a talker's voice and location. Dashed line depicts chance level, asterisks indicate significance: \* =  $p < 0.05$ ; \*\* =  $p < 0.01$ ; \*\*\* =  $p < 0.001$ .

393 **Attention to a localized talker enhances temporal coherence between voice-sensitive and**  
394 **location-sensitive sites.**

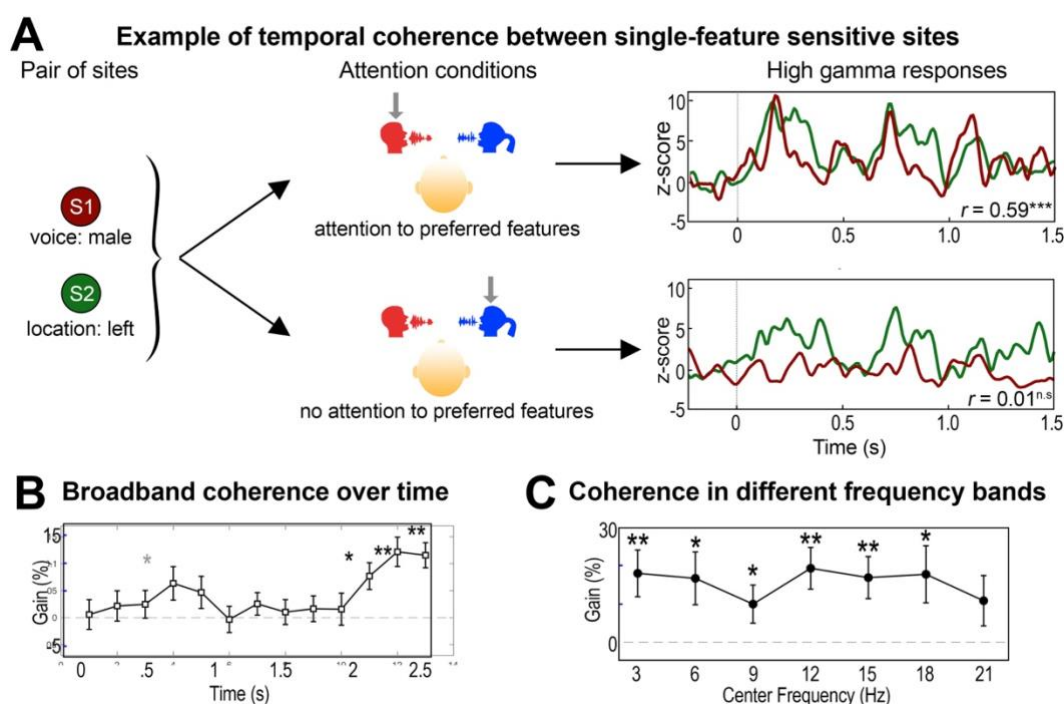
395 We showed that joint population coding results in simultaneous encoding of an attended talker's voice  
396 and location in spatial multi-talker scenes. However, other mechanisms may also contribute to linking  
397 an attended talker's voice and location features in spatial multi-talker scenes. In particular, it has been  
398 proposed that different sound features are bound together through synchronization of the slow  
399 fluctuations in neural responses of feature sensitive cortical sites, that is, temporal coherence<sup>15</sup>. Here, we  
400 evaluated to what extent such temporal coherence contributed to linking the attended talker's voice and  
401 location. First, we computed temporal coherence between the high-gamma envelope of pairs of neural  
402 sites consisting of one site sensitive to a single talker's voice and one site sensitive only to a single talker's  
403 location (Figure 6 A). For each voice-location site pair ( $n = 57$ ), we quantified temporal coherence of the  
404 high gamma envelope at frequencies between 2-22 Hz using the coherency coefficient. The coherency  
405 coefficient is the frequency-domain mathematical equivalent of the cross-correlation function in the  
406 time-domain<sup>38</sup> (within-subjects analysis; Methods). As shown in Figure 6 A, we evaluated the hypothesis  
407 that attention selectively enhances temporal coherence between the voice and location site in each pair  
408 by contrasting temporal coherence in different attention conditions. That is, we examined whether  
409 temporal coherence increased when attention was directed towards a localized talker which matched  
410 the pair's preferred features (condition 'preferred features attended') in comparison to when attention  
411 was directed to a localized talker which was orthogonal to the pair's preferred features (condition  
412 'preferred features unattended'). For example, for a pair of sites consisting of a voice sensitive site tuned  
413 to the *female talker* and a location sensitive site tuned to the *right*, we hypothesize that temporal  
414 coherence increases when attention is directed to a female talker on the right in comparison to when  
415 attention is directed to a male talker on the left and the preferred features are therefore unattended  
416 (Figure 6 A). We quantify such attentional modulation of temporal coherence as the coherence gain:

417 
$$AM_{coh} = \frac{coh_{xy}(\omega)_{attended} - coh_{xy}(\omega)_{unattended}}{coh_{xy}(\omega)_{unattended}}$$

418 Here, we define ‘attended’ as the condition in which attention is directed towards the preferred features  
 419 and the ‘unattended condition’ as the condition in which the preferred features are unattended (i.e.,  
 420 because attention is directed towards orthogonal features).

421 First, we evaluated the development of attentional modulation of temporal coherence over time by  
 422 calculating broadband temporal coherence (i.e. across all frequencies between 2 and 22 Hz) in shifting  
 423 windows of 1,000 ms with 50 % overlap. Figure 6 B shows a weak but not statistically significant  
 424 attentional enhancement of temporal coherence immediately post sound onset (t-test,  $p = 0.074$ , FDR  
 425 corrected) and a strong and robust gain in temporal coherence starting at approximately 2 s post  
 426 stimulus-onset. Furthermore, we examined whether the observed attentional enhancement in the late  
 427 response was generic to the range of frequencies tested here (2-22Hz) or whether it was frequency  
 428 specific. We therefore repeated the analysis on the late response which showed the most robust  
 429 attentional temporal coherence gain (i.e., from 1.75 s until 3.25 s post sound onset) in narrow frequency  
 430 bins of 3 Hz (central frequencies [CF]: 3, 6, 9, 12, 15, 18, 21 Hz). We found that in this time window,  
 431 attentional enhancement of temporal coherence was generic for frequencies < 22 Hz (Figure 6 C).

432 Taken together, these results demonstrate that in spatial multi-talker scenes, attention selectively  
 433 enhanced temporal coherence between sites sensitive to a single talker’s voice and sites sensitive to a  
 434 single talker’s location. Moreover, we showed that this attentional enhancement builds up over time. In  
 435 sum, temporal coherence is a plausible binding mechanism for linking voice and location encoding by  
 436 single-feature sensitive sites in order to form a complete auditory object in complex, multi-source  
 437 auditory scenes.



**Figure 6: Linking an attended talker's voice and location through temporal coherence.** (A) Schematic example of attentional modulation of temporal coherence of a pair of neural sites consisting of a single-feature voice sensitive site (red) and a single-feature location sensitive site (green). (B) Development of broadband temporal coherence gain over time. Error bars reflect standard error of the mean (SEM). Asterisks indicate a significant attentional enhancement of coherence. \*\* =  $p < 0.01$ , \* =  $p < 0.05$ . (C) Temporal coherence gain per narrowband frequency bin.

## 438 **DISCUSSION**

439 In daily-life situations, listeners flexibly extract relevant information from cluttered and dynamic auditory  
440 scenes to form multi-dimensional auditory objects such as a 'localized talker'. While auditory object  
441 formation is critically dependent on the integration of different feature dimensions (e.g. location and  
442 voice), it is presently not clear how such different sound attributes are linked by the brain. Here, we  
443 utilized the unique spatiotemporal resolution of invasive intracranial measurements in neurosurgical  
444 patients to gain insight into the neural mechanisms linking voice and location sound features in real-life  
445 listening scenes consisting of a single talker or two spatially separated talkers.

446 We found that cortical responses varied from dual-feature sensitivity to a talker's voice and location, to  
447 single-feature sensitivity to a talker's voice or location only. Further, population responses of both dual-  
448 feature sensitive and single-feature sensitive sites, simultaneously encoded an attended talker's voice  
449 and location features. Our findings thus indicate that cortical representations of a multi-dimensional  
450 localized talker are derived from joint encoding in distributed population response patterns rather than  
451 separate voice and location encoding in dual processing streams within delineated anatomical  
452 regions<sup>6,7,32</sup>. Furthermore, our data indicate that attention enhances temporal coherence between voice  
453 and location sensitive sites, thereby providing an additional mechanism for linking the representations  
454 of both features. These results provide important new insights into the emergence of multi-dimensional  
455 auditory objects<sup>2,4</sup> in auditory cortex during active, goal-oriented listening in real-life listening scenes.

### 456 **Active task design and naturalistic stimuli reveal distributed voice and location encoding**

457 Our data showed that the sensitivity of local cortical sites for a talker's voice and location features can  
458 be explained by the underlying spectrotemporal tuning properties. These results align with prior research  
459 attributing speaker sensitivity to spectrotemporal tuning properties<sup>21</sup> and fast temporal processing to  
460 the posterior-dorsal regions of human auditory cortex<sup>25</sup> which tend to show strong spatial sensitivity<sup>25</sup>.  
461 Additionally, our results highlight that voice and location encoding during active listening occurs in  
462 distributed networks that span the entire auditory cortex rather than within clearly delineated,  
463 functionally specialized cortical regions. Moreover, linking local responses to population encoding  
464 showed that sites which are characterized by functionally specialized local responses (for example, voice

465 sensitive sites), nevertheless encode information about both voice and location in their population  
466 responses.

467 Further, the distributed networks of voice and location sensitivity that we observed at the level of  
468 individual cortical sites is in agreement with a recent study which demonstrated that acoustic and  
469 phonetic processing in auditory cortex are based on distributed, parallel processing rather than serial  
470 processing<sup>12</sup>. This indicates that distributed processing may be a general characteristic of auditory  
471 encoding and speech encoding specifically<sup>10</sup>. Moreover, the occurrence of distributed voice and location  
472 representations as observed in the present study conceivably ensures sufficient flexibility to  
473 accommodate sound encoding in changing acoustic environments and with changing behavioral  
474 goals<sup>39</sup>.

475 Further, we showed that sensitivity to a talker's location features is similar across sites that are at lower  
476 stages of the hierarchy and sites that are at higher stages of the hierarchy. These results deviate from  
477 the view that spatial sensitivity emerges only in higher-order regions belonging to the functionally  
478 specialized location stream<sup>7,40</sup>. Instead, our findings are in agreement with more recent studies with  
479 active task designs which demonstrated that neural location sensitivity in early processing stages (i.e.  
480 primary auditory cortex) is more pronounced during active, goal-oriented localization<sup>8,9</sup>. Taken together,  
481 our results emphasize that experiment designs comprising active tasks and naturalistic stimuli are crucial  
482 to uncover representational mechanisms related to goal-oriented behavior in complex auditory scenes.

### 483 **Pre-attentive and attentive linking of voice and location to form complete auditory objects**

484 Whether attention is required for auditory object formation remains a matter of debate<sup>2,15</sup>. Some have  
485 argued that auditory streams are formed pre-attentively, for example by the activation of separate  
486 populations of neurons<sup>23</sup>. Others have posited that attention is required to bind together the various  
487 attributes of the attended object<sup>15</sup>. Our data showed that a subset of local cortical sites exhibited  
488 sensitivity to both voice and location features, similar to prior findings of multi-feature sensitivity in ferret  
489 auditory cortex<sup>13</sup>. Moreover, we showed that the population response patterns of these sites gave rise  
490 to representations of the multi-dimensional object, that is, the localized talker. It is therefore conceivable  
491 that the activation of these populations contributes to stream formation in the spatial multi-talker scenes  
492 utilized here. However, to what extent this mechanism is pre-attentive requires further investigation.

493 Our results also revealed top-down attentional modulation of feature binding. That is, a subset of sites  
494 showed single-feature sensitive responses to either voice or location features, in agreement with the  
495 'feature analysis'-stage of the temporal coherence framework. According to this framework, distinct  
496 neural populations generate representations of various sound properties<sup>15,41,42</sup>. Here, we found that  
497 attention selectively enhanced temporal coherence between relevant single-feature voice and location

498 sensitive sites. This result is consistent with accumulating evidence<sup>16,18</sup> supporting temporal coherence  
499 as a potential mechanism for grouping of perceptual features. An open question is where in the cortex  
500 the read-out of such temporally coherent input takes place.

501 Taken together, our data indicate that linking of a talker's voice and location features in spatial multi-  
502 talker scenes emerges from a mixture of (potentially pre-attentive) activation of dual-feature sensitive  
503 neural sites, population coding and attentional modulation of temporal coherence between voice and  
504 location sensitive sites.

### 505 **A continuum of attentional modulations of voice and location encoding**

506 In agreement with prior work<sup>21,24</sup>, our results show that attending to a talker's voice and location elicited  
507 weak attentional response gain control. Further, attention dynamically changed spectrotemporal tuning  
508 properties of late-response cortical sites, resulting in contrast matched filtering shape changes that  
509 enhanced local selective tracking of the attended talker's speech. These results connect prior work in  
510 animals which showed that task performance and attention changed spectrotemporal tuning in auditory  
511 cortex<sup>36,37</sup> to attended speech encoding in complex scenes in human auditory cortex. Moreover, these  
512 results extend findings from prior neural measurements in human auditory cortex, which showed that  
513 contextual information elicited adaptive STRF tuning to boost perception of degraded speech<sup>35</sup>.

514 Further, an attended talker's location features elicited comparable attentional gain control in early- and  
515 late-response sites, suggesting that attention affects spatial processing at low-level as well as higher-  
516 order processing stages. More research is needed to establish whether these spatial attention effects in  
517 low-level cortical regions emerge from feedback projections originating in higher-order regions<sup>39</sup>.  
518 Additional work is also needed to evaluate whether attending a talker's location in multi-talker scenes  
519 affects spatial tuning. That is, studies using single-source experiment designs with an active listening  
520 task reported sharpening of spatial tuning in primary auditory regions<sup>8,9</sup> and it is likely that similar effects  
521 take place in multi-talker scenes to support segregating background from foreground. However, to  
522 assess this hypothesis, an experiment design with more fine-grained sampling of azimuth locations is  
523 required to elucidate attentional modulation of spatial receptive fields.

### 524 **Conclusion and outlook**

525 Our results point to distributed and joint voice and location encoding across auditory cortex during  
526 active, goal-directed behavior. These findings support the view that object formation and attentional  
527 selection emerge gradually and in a distributed manner from the auditory hierarchy, rather than at one  
528 specific site or region in auditory cortex<sup>4</sup>. Such a distributed code flexibly accommodates rapid changes  
529 in the (acoustic) environment as well as changing behavioral goals. Crucially, the present findings  
530 demonstrate the need for real-life, complex stimuli and experimental designs including active behavioral

531 tasks to understand cortical processing of multi-dimensional auditory objects. Future studies including  
532 stimuli spanning a larger and more fine-grained range of talkers, locations and other sound features can  
533 further unravel local cortical tuning properties as well as population representations of multi-  
534 dimensional auditory objects. Finally, complementing sEEG measurements with high-density intracranial  
535 measurements (e.g. high-density electrocorticography [ECoG], e.g. <sup>43</sup>) are critical to refine cortical maps  
536 of local feature sensitivity, to tease apart fine-grained population representations within and across  
537 macro-anatomical regions, and to further our insights into feature binding through temporal coherence.

538 **REFERENCES**

- 539 1. Bregman, A. S. *Auditory Scene Analysis: The Perceptual Organization of Sound*. (MIT press,  
540 1994).
- 541 2. Bizley, J. K. & Cohen, Y. E. The what, where and how of auditory-object perception. *Nat Rev*  
542 *Neurosci* **14**, 693–707 (2013).
- 543 3. Darwin, C. J. Auditory grouping. *Trends Cogn Sci* **1**, 327–333 (1997).
- 544 4. Shinn-Cunningham, B., Best, V. & Lee, A. K. C. Auditory object formation and selection. in *The*  
545 *auditory system at the cocktail party* 7–40 (Springer, 2017).
- 546 5. Bizley, J. K., Walker, K. M. M., Nodal, F. R., King, A. J. & Schnupp, J. W. H. Auditory cortex  
547 represents both pitch judgments and the corresponding acoustic cues. *Current Biology* **23**,  
548 620–625 (2013).
- 549 6. Tian, B., Reser, D., Durham, A., Kustov, A. & Rauschecker, J. P. Functional specialization in  
550 rhesus monkey auditory cortex. *Science (1979)* **292**, 290–293 (2001).
- 551 7. Rauschecker, J. P. & Tian, B. Mechanisms and streams for processing of “what” and “where”  
552 in auditory cortex. *Proceedings of the National Academy of Sciences* **97**, 11800–11806 (2000).
- 553 8. Lee, C.-C. & Middlebrooks, J. C. Auditory cortex spatial sensitivity sharpens during task  
554 performance. *Nat Neurosci* **14**, 108–114 (2011).
- 555 9. van der Heijden, K., Rauschecker, J. P., Formisano, E., Valente, G. & de Gelder, B. Active sound  
556 localization sharpens spatial tuning in human primary auditory cortex. *Journal of Neuroscience*  
557 **38**, 8574–8587 (2018).
- 558 10. Bhaya-Grossman, I. & Chang, E. F. Speech computations of the human superior temporal  
559 gyrus. *Annu Rev Psychol* **73**, 79–102 (2022).
- 560 11. Yi, H. G., Leonard, M. K. & Chang, E. F. The encoding of speech sounds in the superior  
561 temporal gyrus. *Neuron* **102**, 1096–1110 (2019).
- 562 12. Hamilton, L. S., Oganian, Y., Hall, J. & Chang, E. F. Parallel and distributed encoding of speech  
563 across human auditory cortex. *Cell* **184**, 4626–4639 (2021).
- 564 13. Bizley, J. K., Walker, K. M. M., Silverman, B. W., King, A. J. & Schnupp, J. W. H. Interdependent  
565 encoding of pitch, timbre, and spatial location in auditory cortex. *Journal of neuroscience* **29**,  
566 2064–2075 (2009).
- 567 14. King, A. J., Teki, S. & Willmore, B. D. B. Recent advances in understanding the auditory cortex.  
568 *F1000Res* **7**, (2018).
- 569 15. Shamma, S. A., Elhilali, M. & Micheyl, C. Temporal coherence and attention in auditory scene  
570 analysis. *Trends Neurosci* **34**, 114–123 (2011).
- 571 16. Lu, K. *et al.* Temporal coherence structure rapidly shapes neuronal interactions. *Nat Commun*  
572 **8**, 1–12 (2017).
- 573 17. Elhilali, M., Ma, L., Micheyl, C., Oxenham, A. J. & Shamma, S. A. Temporal coherence in the  
574 perceptual organization and cortical representation of auditory scenes. *Neuron* **61**, 317–329  
575 (2009).

- 576 18. O’Sullivan, J. A., Shamma, S. A. & Lalor, E. C. Evidence for neural computations of temporal  
577 coherence in an auditory scene and their enhancement during active listening. *Journal of*  
578 *Neuroscience* **35**, 7256–7263 (2015).
- 579 19. King, A. J. & Walker, K. M. M. Listening in complex acoustic scenes. *Curr Opin Physiol* **18**, 63–  
580 72 (2020).
- 581 20. Fritz, J. B., Elhilali, M., David, S. V & Shamma, S. A. Auditory attention—focusing the  
582 searchlight on sound. *Curr Opin Neurobiol* **17**, 437–455 (2007).
- 583 21. O’Sullivan, J. *et al.* Hierarchical encoding of attended auditory objects in multi-talker speech  
584 perception. *Neuron* **104**, 1195–1209 (2019).
- 585 22. Mesgarani, N. & Chang, E. F. Selective cortical representation of attended speaker in multi-  
586 talker speech perception. *Nature* **485**, 233–236 (2012).
- 587 23. Sussman, E. S., Horváth, J., Winkler, I. & Orr, M. The role of attention in the formation of  
588 auditory streams. *Percept Psychophys* **69**, 136–152 (2007).
- 589 24. Patel, P. *et al.* Interaction of bottom-up and top-down neural mechanisms in spatial multi-  
590 talker speech perception. *Current Biology* **32**, 3971–3986 (2022).
- 591 25. Derey, K., Valente, G., de Gelder, B. & Formisano, E. Opponent coding of sound location  
592 (azimuth) in planum temporale is robust to sound-level variations. *Cerebral Cortex* **26**, 450–  
593 464 (2016).
- 594 26. Allen, E. J., Burton, P. C., Olman, C. A. & Oxenham, A. J. Representations of pitch and timbre  
595 variation in human auditory cortex. *Journal of neuroscience* **37**, 1284–1293 (2017).
- 596 27. Steinschneider, M., Nourski, K. V, Rhone, A. E., Kawasaki, H. & Oya, H. Differential activation  
597 of human core, non-core and auditory-related cortex during speech categorization tasks as  
598 revealed by intracranial recordings. *Front Neurosci* **8**, 103289 (2014).
- 599 28. Schnupp, J., Nelken, I. & King, A. *Auditory Neuroscience: Making Sense of Sound*. (MIT press,  
600 2011).
- 601 29. Risoud, M. *et al.* Sound source localization. *Eur Ann Otorhinolaryngol Head Neck Dis* **135**, 259–  
602 264 (2018).
- 603 30. Moerel, M., De Martino, F. & Formisano, E. An anatomical and functional topography of  
604 human auditory cortical areas. *Front Neurosci* **8**, 225 (2014).
- 605 31. Nourski, K. V *et al.* Functional organization of human auditory cortex: investigation of  
606 response latencies through direct recordings. *Neuroimage* **101**, 598–609 (2014).
- 607 32. Rauschecker, J. P. & Scott, S. K. Maps and streams in the auditory cortex: nonhuman primates  
608 illuminate human speech processing. *Nat Neurosci* **12**, 718–724 (2009).
- 609 33. Rifkin, R., Yeo, G. & Poggio, T. Regularized least-squares classification. *Nato Science Series Sub*  
610 *Series III Computer and Systems Sciences* **190**, 131–154 (2003).
- 611 34. Golumbic, E. M. Z. *et al.* Mechanisms underlying selective neuronal tracking of attended  
612 speech at a “cocktail party”. *Neuron* **77**, 980–991 (2013).
- 613 35. Holdgraf, C. R. *et al.* Rapid tuning shifts in human auditory cortex enhance speech  
614 intelligibility. *Nat Commun* **7**, 13654 (2016).



- 615 36. Fritz, J., Shamma, S., Elhilali, M. & Klein, D. Rapid task-related plasticity of spectrotemporal  
616 receptive fields in primary auditory cortex. *Nat Neurosci* **6**, 1216–1223 (2003).
- 617 37. Fritz, J. B., Elhilali, M., David, S. V & Shamma, S. A. Does attention play a role in dynamic  
618 receptive field adaptation to changing acoustic salience in A1? *Hear Res* **229**, 186–203 (2007).
- 619 38. Bastos, A. M. & Schoffelen, J.-M. A tutorial review of functional connectivity analysis methods  
620 and their interpretational pitfalls. *Front Syst Neurosci* **9**, 175 (2016).
- 621 39. van der Heijden, K., Rauschecker, J. P., de Gelder, B. & Formisano, E. Cortical mechanisms of  
622 spatial hearing. *Nat Rev Neurosci* **20**, 609–623 (2019).
- 623 40. Alain, C., Arnott, S. R., Hevenor, S., Graham, S. & Grady, C. L. “What” and “where” in the  
624 human auditory system. *Proceedings of the national academy of sciences* **98**, 12301–12306  
625 (2001).
- 626 41. Fishman, Y. I., Arezzo, J. C. & Steinschneider, M. Auditory stream segregation in monkey  
627 auditory cortex: effects of frequency separation, presentation rate, and tone duration. *J*  
628 *Acoust Soc Am* **116**, 1656–1670 (2004).
- 629 42. Fishman, Y. I., Reser, D. H., Arezzo, J. C. & Steinschneider, M. Neural correlates of auditory  
630 stream segregation in primary auditory cortex of the awake monkey. *Hear Res* **151**, 167–187  
631 (2001).
- 632 43. Ramsey, N. F. *et al.* Decoding spoken phonemes from sensorimotor cortex with high-density  
633 ECoG grids. *Neuroimage* **180**, 301–311 (2018).
- 634 44. Kayser, C., Petkov, C. I. & Logothetis, N. K. Tuning to sound frequency in auditory field  
635 potentials. *J Neurophysiol* **98**, 1806–1809 (2007).
- 636 45. Nir, Y. *et al.* Coupling between neuronal firing rate, gamma LFP, and BOLD fMRI is related to  
637 interneuronal correlations. *Current biology* **17**, 1275–1285 (2007).
- 638 46. Yang, X., Wang, K. & Shamma, S. A. Auditory representations of acoustic signals. *IEEE Trans Inf*  
639 *Theory* **38**, 824–839 (1992).
- 640 47. Chi, T., Ru, P. & Shamma, S. A. Multiresolution spectrotemporal analysis of complex sounds. *J*  
641 *Acoust Soc Am* **118**, 887–906 (2005).
- 642 48. Theunissen, F. E. *et al.* Estimating spatio-temporal receptive fields of auditory and visual  
643 neurons from their responses to natural stimuli. *Network: Computation in Neural Systems* **12**,  
644 289 (2001).
- 645 49. Santoro, R. *et al.* Encoding of natural sounds at multiple spectral and temporal resolutions in  
646 the human auditory cortex. *PLoS Comput Biol* **10**, e1003412 (2014).
- 647 50. Cooke, M. A glimpsing model of speech perception in noise. *J Acoust Soc Am* **119**, 1562–1573  
648 (2006).
- 649
- 650

## 651 **METHODS**

### 652 **Preprocessing**

653 A detailed description of preprocessing of the neural data can be found in <sup>24</sup>. In short, data preprocessing  
654 included montaging to a common average reference, noise removal, extraction of the high gamma  
655 envelope (70 – 150 Hz) using the Hilbert transform. The high-gamma envelope is thought to reflect  
656 neuronal population activity<sup>44,45</sup>. Finally, neural responses were down sampled to 100 Hz and z-scored  
657 across single speaker blocks and across multi-talker blocks (i.e. calculated over both male and female  
658 trials, but separately for single- and multi-source blocks).

### 659 **Speech responsive electrodes**

660 To assess which electrodes exhibited a robust response to speech streams, we computed for each  
661 electrode the mean baseline response as the average of the high-gamma envelope during 0.5 seconds  
662 preceding stimulus onset, and the mean speech onset response as the average of the high-gamma  
663 envelope in the 0.5 seconds following stimulus onset. To test for a statistically significant auditory  
664 response, we performed a paired samples t-test for each electrode and applied FDR correction across  
665 electrodes to correct for multiple comparisons. Only electrodes that exhibited a robust auditory response  
666 at  $q < 0.05$  were included in the remainder of the analysis.

### 667 **Estimating spectrotemporal receptive fields (STRFs) and response latency**

668 First, we computed a cortical spectrogram representation of each sound scene using a model of early  
669 cochlear processing and mid-brain auditory processing (NSL toolbox<sup>1</sup>). We modeled cochlear processing  
670 using a filter bank of 128 constant-Q filters that were spaced equally on a logarithmic axis ranging from  
671 center frequency (CF) = 270 Hz to CF = 7,246 Hz. Next, we modeled auditory midbrain processing by  
672 taking the derivative along the frequency axis, performing half-wave rectification and applying short-  
673 term temporal integration<sup>46,47</sup>. This approach accounted for the enhanced frequency selectivity as a  
674 consequence of lateral inhibition, as well as reduced phase locking, observed after midbrain processing.  
675 Cortical spectrograms were computed based on monaural stimulus waveforms (i.e. independent of  
676 sound location). The resulting spectrograms had a sampling frequency of 100 Hz and were down  
677 sampled to 50 channels to reduce the number of parameters.

678 We then estimated the spectrotemporal receptive field (STRF) by linearly mapping the cortical  
679 spectrogram to the evoked response using the STRFlab MATLAB Toolbox<sup>48</sup> (<http://strflab.berkeley.edu>).  
680 For each electrode, we used the past 300 ms of a stimulus to predict the neural response at every time  
681 point using normalized reverse correlation. To prevent overfitting, we used a five-fold cross-validation

---

<sup>1</sup> Available from <http://nsl.isr.umd.edu/index.html>

682 procedure. We optimized sparsity and regularization parameters by maximizing the correlation between  
683 actual and predicted responses. Using the resulting STRFs, we defined the response latency for each  
684 electrode as the time point corresponding to the peak energy in the STRF.

### 685 **Decoding a single talker's voice and location features**

686 We trained a four-class classifier on population neural response patterns to jointly decode a talker's  
687 voice and location features. The four classes corresponded to 'male talker, left', 'male talker, right',  
688 'female talker, left' and 'female talker, right'. We used frame-by-frame, regularized least-squares (RLS)  
689 classification<sup>22,33</sup> which produced for each time frame a linear weighted sum of the population of neural  
690 responses for each class<sup>22</sup>. We trained and tested classifiers on the sustained responses only (i.e.,  
691 excluding response onset effects from 0 to 500 ms). The class with the highest average classifier output  
692 over all frames in the trial was taken as the predicted class.

693 We trained the classifier in a leave-two-trials-out cross-validation procedure on the single-talker data  
694 (corresponding to 25 folds). We computed classification accuracy as the average accuracy across the 25  
695 folds. Further, to evaluate the statistical significance of classification accuracies, we performed a  
696 permutation analysis in which we randomly permuted the class labels and repeated the complete 25-  
697 fold cross-validation procedure. We iterated this process 2,000 times to create a null distribution of  
698 classification accuracy. Next, we tested whether the observed classification accuracy exceeds the 95<sup>th</sup>  
699 percentile of the null distribution of permuted accuracies (one sided test). We computed  $p$  as the  
700 proportion of permuted accuracies that was equal to or larger than the observed accuracy.

701 Finally, we calculated marginal accuracies for the voice and location feature dimensions by labelling  
702 accuracy based on a single feature dimension only, ignoring the other feature dimension. For example,  
703 to quantify the marginal accuracy for voice features, we calculated the percentage of the trials for which  
704 the correct voice class was predicted (i.e. female or male talker), ignoring the predicted location class  
705 (i.e. left or right). We computed the marginal accuracy also as the average across the 25 folds and used  
706 the permutation procedure described above to assess the statistical significance of the marginal  
707 accuracies.

### 708 **Attentional-driven response gains in multi-talker scenes**

709 For each electrode, we quantified the strength and direction of attentional modulations of cortical  
710 responses in the multi-source scenes evoked either by attending to a talker's voice features or by  
711 attending to a talker's location features using Cohen's  $d$ . That is, similar to the quantification of single-  
712 talker feature sensitivity described above, we computed the mean response for each trial in the multi-  
713 source condition as the mean from 0.5 s post sound onset to 1.5 s post sound onset. Then, to test for  
714 attention-driven response gains for a talker's voice features, we computed the effect size for the

715 difference between the mean responses to all ‘attend male’ and ‘attend female’ trials, irrespective of the  
716 attended location of the trials ( $n = 50$  each). To test for attention-driven response gains for a talker’s  
717 location features, we computed the effect size for the difference between the mean responses to all  
718 ‘attend left’ and ‘attend right’ trials, irrespective of the attended talker of the trials ( $n = 50$  each).

### 719 **Decoding an attended talker’s voice and location features in multi-talker scenes**

720 To decode an attended talker’s voice and location features in spatial multi-talker scenes, we trained the  
721 four-class classifier on the multi-talker data using a similar procedure as described above. In multi-talker  
722 scenes, class labels consisted of ‘attended male talker, left’, ‘attended male talker, right’, attended female  
723 talker, left’ and ‘attended female talker, right’. We also assessed statistical significance using a  
724 permutation procedure similar to the permutation procedure for single talker data.

### 725 **Quantifying temporal coherence**

726 We assessed temporal coherence in slow fluctuations in stimulus evoked responses between pairs of  
727 voice sensitive and location sensitive sites. This analysis was performed on a within subject level. Five  
728 subjects contained multiple voice-location pairs of and were therefore included in the analysis. Because  
729 the high-gamma envelope is considered a signature of neural population responses<sup>44,45</sup>, we computed  
730 temporal coherence on the high-gamma envelope. Further, we quantified temporal coherence using the  
731 coherency coefficient, which is the mathematical equivalent in the frequency domain of the cross-  
732 correlation function in the time domain<sup>38</sup>. Specifically, the coherence coefficient is the normalized  
733 average cross-power spectral density between signals  $x$  and  $y$  across trials at frequency  $\omega$  computed  
734 as<sup>38</sup>:

$$735 \quad coh_{xy}(\omega) = \frac{\left| \frac{1}{n} \sum_{k=1}^n A_x(\omega, k) A_y(\omega, k) e^{i(\phi_x(\omega, k) - \phi_y)} \right|}{\sqrt{\left( \left( \frac{1}{n} \sum_{k=1}^n A_x^2(\omega, k) \right) \left( \frac{1}{n} \sum_{k=1}^n A_y^2(\omega, k) \right) \right)}}$$

736 Here, we computed broadband temporal coherence over a frequency range of 2 – 22 Hz to map the  
737 development of attentional enhancement of temporal coherence over time, correspond to the range of  
738 slow fluctuations in which temporal coherence for feature binding is hypothesized to occur (i.e. 50 ms  
739 to 500 ms<sup>15</sup>). Furthermore, we computed narrowband temporal coherency for eight frequency bands  
740 with center frequencies 3, 6, 9, 12, 15, 18, and 21 Hz (bandwidth = 3 Hz) to examine the effect of attention  
741 on temporal coherence for specific frequencies.

742

743 **ACKNOWLEDGEMENTS**

744 We thank Menoua Keshishian for sharing his expertise on STRF analysis. This study was supported by  
745 National Institute on Deafness and other Communication Disorders grant R01DC014279 (NM) and grant  
746 R01DC018805 (NM). This project has also received funding from the European Union's Horizon 2020  
747 Research and Innovation Program under the Marie Skłodowska-Curie grant agreement No 898134 (KH)  
748 and from the NWO Talent Program under the Veni grant agreement VI.Veni.202.184 (KH).