

1 *The following article has been submitted to The Journal of the Acoustical Society of America. After*  
2 *it is published, it will be found at <http://asa.scitation.org/journal/jas>*

3

4

## **Perceptual grouping in the cocktail party:**

5

## **contributions of voice-feature continuity**

6

7 Jens Kreitewolf<sup>f\*1,2</sup>, Samuel R. Mathias<sup>3</sup>, Régis Trapeau<sup>1,4</sup>, Jonas Obleser<sup>2</sup>, Marc Schönwiesner<sup>1,5</sup>

8

9 *<sup>1</sup>International Laboratory for Brain, Music and Sound Research (BRAMS),*

10 *Department of Psychology, Université de Montréal, Canada; <sup>2</sup>Department of Psychology,*

11 *University of Lübeck, Germany; <sup>3</sup>Neurocognition, Neurocomputation and Neurogenetics (n3)*

12 *Division, Yale University School of Medicine, USA; <sup>4</sup>Institut de Neurosciences de la Timone UMR*

13 *7289, Centre National de la Recherche Scientifique and Aix-Marseille Université, France; <sup>5</sup>Institute*

14 *of Biology, Leipzig University, Germany*

15

16

17

18

19 *\*Correspondence should be addressed to:*

20 Jens Kreitewolf

21 Department of Psychology, University of Lübeck, Maria-Goeppert-Str. 9a, 23562 Lübeck,

22 Germany; E-mail: [jens.kreitewolf@uni-luebeck.de](mailto:jens.kreitewolf@uni-luebeck.de)

23 **Abstract**

24 Cocktail parties pose a difficult yet solvable problem for the auditory system. Previous work  
25 has shown that the cocktail-party problem is considerably easier when all sounds in the target  
26 stream are spoken by the same talker (the *voice-continuity benefit*). The present study  
27 investigated the contributions of two of the most salient voice features — glottal-pulse rate  
28 (GPR) and vocal-tract length (VTL) — to the voice-continuity benefit. Twenty young, normal-  
29 hearing listeners participated in two experiments. On each trial, listeners heard concurrent  
30 sequences of spoken digits from three different spatial locations and reported the digits  
31 coming from a target location. Critically, across conditions, GPR and VTL either remained  
32 constant or varied across target digits. Additionally, across experiments, the target location  
33 either remained constant (Experiment 1) or varied (Experiment 2) within a trial. In Experiment  
34 1, listeners benefited from continuity in either voice feature, but VTL continuity was more  
35 helpful than GPR continuity. In Experiment 2, spatial discontinuity greatly hindered listeners'  
36 abilities to exploit continuity in GPR and VTL. The present results suggest that selective  
37 attention benefits from continuity in target voice features, and that VTL and GPR play  
38 different roles for perceptual grouping and stream segregation in the cocktail party.

39

40 **Keywords:** cocktail party; perceptual grouping; glottal-pulse rate; vocal-tract length

## 41 **I. Introduction**

42 In everyday life, sounds rarely occur in isolation; instead, most of time, the auditory scene  
43 comprises a multitude of sounds heard at once. Consequently, the auditory signal that  
44 reaches the listener's ears is usually a mixture of sounds elicited by various sources. These  
45 situations are often referred to as *cocktail parties* (Cherry, 1953) and pose a difficult  
46 conceptual problem for the listener: to ensure comprehension of target speech, listeners  
47 need to attend to the target voice while at the same time ignoring other irrelevant sounds.

48 Previous work has shown that the cocktail-party problem is made considerably easier  
49 when all target sounds are spoken by the same talker (Best et al., 2008; Bressler et al., 2014;  
50 Kitterick et al., 2010; Larson and Lee, 2013). In the following, we refer to this phenomenon as  
51 the *voice-continuity benefit*. The voice-continuity benefit occurs because speech sounds from  
52 a single talker are all similar in terms of certain acoustic features, which makes it easier to  
53 perceptually group together these sounds than speech sounds produced by different talkers.  
54 Importantly, previous studies demonstrating the voice-continuity benefit all used natural  
55 speech. It is therefore unclear precisely which features common to speech sounds produced  
56 by the same talker contribute to the voice-continuity benefit.

57 A separate line of research has investigated which features are important for  
58 distinguishing different talkers and recognizing familiar ones (reviewed by Mathias and von  
59 Kriegstein, 2014). This work has shown that two of the most salient features are glottal-pulse  
60 rate (GPR) and vocal-tract length (VTL). GPR is the oscillation rate of the vocal folds; it  
61 determines the fundamental frequency ( $f_0$ ) of a speech sound and is perceived as vocal pitch.  
62 VTL is correlated with a talker's perceived height or body size (e.g., Smith et al., 2005); it

63 determines the spectral envelope of a speech sound and is perceived as an aspect of vocal  
64 timbre. GPR and VTL appear to be the most important cues for rating the similarity of speech  
65 sounds produced by unfamiliar talkers (e.g., Baumann and Belin, 2010; Gaudrain et al., 2009)  
66 and for identifying personally familiar talkers (e.g., Lavner et al., 2000).

67 Previous studies indicate that listeners use GPR and VTL information during cocktail-  
68 party listening. Darwin et al. (2003) presented listeners with two concurrent sentences that  
69 differed in GPR and/or VTL and asked them to report key words from a target sentence. Their  
70 results showed that differences in both GPR and VTL helped listeners to segregate target and  
71 masker sentences, and that differences in both GPR and VTL that were large enough to  
72 simulate a shift in the perception of talker sex helped segregation more than differences in  
73 either GPR or VTL alone. In another study, Vestergaard et al. (2009) showed that, when no  
74 other cues are available for stream segregation, smaller differences in GPR than VTL were  
75 necessary to yield the same performance. Thus, their results suggest that GPR is the more  
76 important cue for stream segregation.

77 Solving the cocktail-party problem requires both *segregation* (separating sounds from  
78 different sources) and *grouping* (binding successive sounds from the same source) (Bregman,  
79 1990). While these previous studies provide evidence for the importance of GPR and VTL for  
80 stream segregation, there is, to date, no direct evidence as to whether GPR and VTL are also  
81 important for perceptual grouping in the cocktail party.

82 The main objective of the present study was therefore to investigate the roles of GPR  
83 and VTL for perceptual grouping by determining their relative contributions to the voice-  
84 continuity benefit. To this end, our experimental manipulations did not concern differences in

85 GPR and VTL across target and masker streams, as in previous studies (Darwin et al., 2003;  
86 Vestergaard et al., 2009), but the continuity of GPR and VTL *within* the target stream.

87         We conducted two experiments with similar designs, involving the same listeners. In  
88 both experiments, listeners heard streams of spoken digits presented simultaneously from  
89 different locations and reported the digits from a target location (Fig. 1A). To explore the  
90 contributions of GPR and VTL to the voice-continuity benefit, we manipulated continuity in  
91 GPR and/or VTL in the target stream (Fig. 1B). This was done by resynthesizing original  
92 recordings of spoken digits using vocoder software (Kawahara et al., 2008). If GPR and VTL are  
93 used for perceptual grouping, listeners should benefit from continuity in these features; that  
94 is, they should report more target digits when GPR and VTL are continuous across target  
95 digits than when they are not. To quantify the benefits from continuity in either GPR, VTL, or  
96 both, we compared the proportions of correctly reported target digits across conditions.  
97 Furthermore, to explore whether continuity in certain voice features helps listeners to “tune  
98 into” the target stream, we compared the probabilities of correctly reporting the current  
99 target conditioned on whether or not the previous target digit was correctly reported  
100 (Bressler et al., 2014).

101         In addition to voice continuity, spatial continuity plays an important role for  
102 perceptual grouping in the cocktail party. Previous studies have shown that performance  
103 deteriorates drastically when listeners are uncertain about the locations of the target sounds  
104 (Best et al., 2008; Brungart and Simpson, 2007; Kidd et al., 2005; Kitterick et al., 2010). In the  
105 present study, we sought to extend these findings by comparing the benefits from voice-

106 feature continuity across two experiments that differed with respect to spatial continuity in  
107 the target stream (Fig. 1A).  
108 Specifically, the comparison between experiments allowed us to investigate whether spatial  
109 discontinuity mediates listeners' abilities to exploit target voice features. Previous work has  
110 shown that listeners only benefit from knowledge about the target voice when the cocktail  
111 party is challenging enough (Kitterick et al., 2010). If spatial discontinuity made the cocktail  
112 party more challenging, one might hypothesize that listeners would attain *greater* benefits  
113 from continuity in GPR and VTL when the target location varies within a trial. On the other  
114 hand, it has been suggested that the temporal coherence across perceptual features,  
115 including pitch, timbre, and location, is crucial for auditory scene analysis (Shamma et al.,  
116 2011). Hence, an alternative possibility is that the lack of spatial continuity would prevent  
117 listeners from fully exploiting continuity in GPR and VTL; if true, we would observe *smaller*  
118 benefits from voice-feature continuity when the target location varies.

## 119 **II. Methods**

### 120 **A. Listeners**

121 Twenty listeners (12 females, 8 males; age: 18-31 years) participated in two experiments. All  
122 listeners were native English speakers and had hearing thresholds of 20 dB hearing level (HL)  
123 or lower at octave-spaced frequencies between 0.125 and 8 kHz. None of them had a history  
124 of hearing disorder or neurological disease. Written informed consent was obtained from all  
125 listeners according to procedures approved by the Research Ethics Committee of the  
126 Université de Montréal. Listeners completed each of two sessions within 1.5 h and were  
127 compensated C\$ 12.50/h for their participation.

## 128 **B. Stimuli**

129 The stimuli were based on the digits one to nine spoken by a native English male talker. Digit  
130 seven was excluded from the stimulus set because it was disyllabic. All other digits were  
131 resynthesized using vocoder software (TANDEM-STRAIGHT; Kawahara et al., 2008) to simulate  
132 nine virtual talkers with different GPRs (79, 150, 285 Hz) and VTLs (8, 12, 18 cm). These values  
133 conform to a stepwise increase of approximately 90 % in GPR and 50 % in VTL. Previous work  
134 has shown that listeners perceive different talker identities at half of these step sizes  
135 (Gaudrain et al., 2009). The loudness of all stimuli was normalized using the Zwicker and Fastl  
136 (1999) model as implemented in the Genesis Loudness Toolbox ([www.genesis.fr](http://www.genesis.fr)). This  
137 procedure also included shifting waveforms in time to ensure that all stimuli were  
138 isochronous.

139 Finally, the stimuli were concatenated into four-digit sequences with an inter-digit  
140 delay of 50 ms using MATLAB (MathWorks, Natick, MA, USA). The digit sequences were  
141 presented through loudspeakers (Orb Audio, New York, NY, USA) via digital-to-analogue  
142 conversion hardware (Tucker-Davis Technologies, Alachua, FL, USA) at 65 dB SPL and with a  
143 sampling rate of 48.828 kHz.

## 144 **C. Apparatus**

145 The study took place in a hemi-anechoic room (2.5 x 5.5 x 2.5 m). Listeners were seated in a  
146 comfortable chair, located in the center of a spherical array of 80 loudspeakers with a  
147 diameter of 1.8 m. Each loudspeaker was equipped with a light-emitting diode (LED).  
148 Listeners were instructed to focus on the central loudspeaker during sound presentation. This  
149 was controlled by a laser pointer and an electromagnetic head-tracking sensor that were

150 attached to the listeners' forehead via a headband. An Optimus Maximus keyboard (Art.  
151 Lebedev Studio, Moscow, Russia) with only the numbers 1 to 9 (excluding 7) lid up on the  
152 number pad was placed on the listeners' lap and served as a response device. Listeners were  
153 instructed to look down at the keyboard to make their responses. Following the listeners'  
154 response, sound presentation only continued once the listeners re-aligned their head with  
155 the central loudspeaker. In case of head misalignment, a 150-Hz tone was played for 200 ms.

## 156 **D. Procedure**

157 The study comprised two sessions conducted on two separate days. In the first session,  
158 listeners were familiarized with the stimuli and the equipment of the main experiments  
159 before performing the two experiments. In the second session, the listeners repeated the two  
160 experiments and finally performed an adaptive procedure that estimated just-noticeable  
161 differences (JNDs) for GPR and VTL.

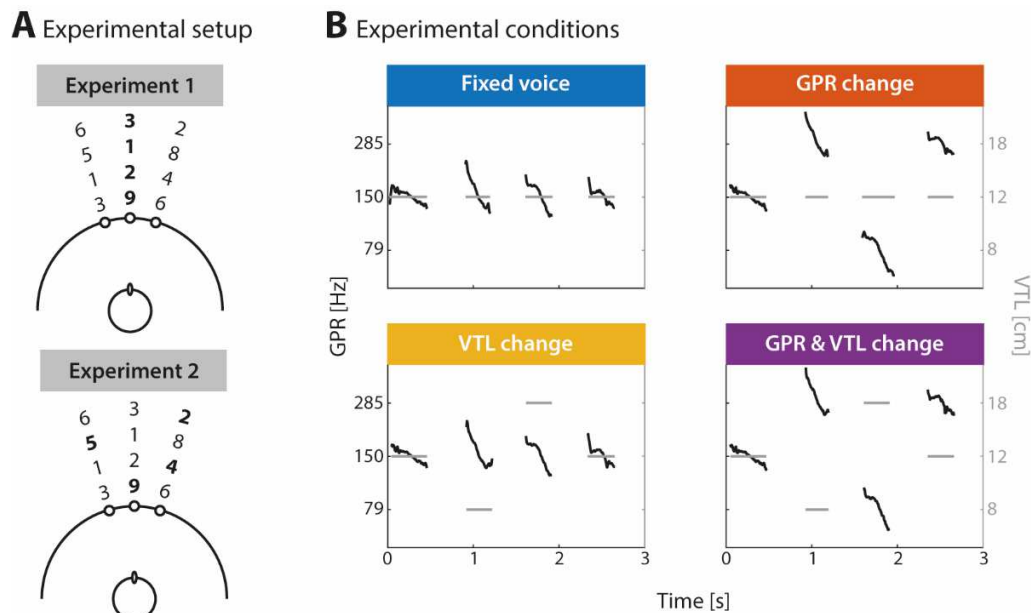
### 162 **1. Experiment 1**

163 On each trial of Experiment 1, listeners heard three competing digit sequences presented  
164 through loudspeakers located at  $-15^\circ$ ,  $0^\circ$ , and  $15^\circ$  on the azimuth (Fig. 1A, upper panel). An  
165 LED affixed to the center of each loudspeaker was illuminated during sound presentation to  
166 indicate the position of the target digit. There was no delay between sound and light onset.  
167 The position of all digits in a target sequence was fixed. The listeners' task was to report the  
168 digits of the target sequence in the order of their presentation. Responses were self-paced  
169 and only allowed after the entire sequence was played. No feedback was given.

170 To investigate the contributions of different voice features to the voice-continuity  
171 benefit, the experiment employed four conditions that differed in terms of (dis-)continuity in



172 GPR and VTL in the target sequence (Fig. 1B): the digits in the target sequence were either  
173 spoken by the same virtual talker (Fixed voice) or virtual talkers whose voices differed in GPR  
174 only (GPR change), in VTL only (VTL change), or in both GPR and VTL (GPR & VTL change). The  
175 order of conditions was pseudo-randomized with the restriction that each condition occurred  
176 within blocks of four trials. The experiment comprised 36 trials per condition in each session.  
177 Figure 1A (upper panel) shows an example trial in Experiment 1. Here, the target sequence  
178 was presented through the central loudspeaker. Throughout the experiment, the target  
179 sequence in each of the four conditions was presented equally often at each of the three  
180 loudspeaker positions. The three concurrently presented digits were always different from  
181 one another and spoken by three different virtual talkers. Furthermore, we ensured that in the  
182 target sequence as well as in each of the two masker sequences, each individual stimulus (i.e.,  
183 each of the eight digits spoken by each of the nine talkers) occurred equally often throughout  
184 the experiment. To familiarize the listeners with the procedure of Experiment 1, they first  
185 conducted a practice block in each of the two sessions. The practice block comprised two  
186 trials of each condition. After completion of Experiment 1, the listeners could take a longer  
187 break before continuing with Experiment 2.



188

189 **Figure 1. (A)** Setup of Experiments 1 (upper panel) and 2 (lower panel). Different four-digit sequences were  
 190 presented simultaneously through three loudspeakers (indicated by full circles on the semi-circle). Time is  
 191 represented as the distance from the loudspeakers. Bold face indicates target digits. In Exp. 1, all digits within a  
 192 target sequence were presented at the same location. In Exp. 2, the target locations varied from digit to digit. All  
 193 three target locations were equiprobable in each condition and each experiment. **(B)** Experimental conditions.  
 194 Digits in the target sequence were either spoken by the same talker (Fixed voice), talkers whose voices differed in  
 195 GPR only (GPR change), in VTL only (VTL change), or in both GPR and VTL (GPR & VTL change). GPR and VTL are  
 196 shown as black  $f_0$  contours and gray bars, respectively.

197

## 198 2. Experiment 2

199 Experiment 2 was designed to investigate the influence of spatial discontinuity on the  
 200 listeners' abilities to group sounds based on voice-feature continuity. In Experiment 1, all  
 201 digits within the target stream were presented at the same location. In Experiment 2, the  
 202 target location varied from digit to digit (Fig. 1A, lower panel). We ensured that there was  
 203 always a change in target location between two consecutive digits and that each of the three  
 204 possible target locations was used at least once per trial. Otherwise, Experiment 2 was  
 205 identical to Experiment 1. Like in Experiment 1, the listeners first conducted a practice block  
 206 in each of the two sessions to familiarize themselves with the experimental procedure. All

207 listeners completed both experiments. However, data from one listener in Experiment 2 were  
208 not recorded in the first session due to technical issues and were dropped in the data analysis.

### 209 **3. Assessment of JNDs**

210 To assess individual listeners' sensitivity to changes in GPR and VTL, we measured just-  
211 noticeable differences (JNDs). For both GPR- and VTL-JNDs, we used a weighted one-up one-  
212 down adaptive procedure that estimates 75 %-correct on the psychometric function  
213 (Kaernbach, 1991). On each trial, two versions of the spoken digit nine were played in  
214 succession from the central loudspeaker (with an inter-stimulus interval of 200 ms).

215 To assess JNDs for GPR, the two versions of the digit nine differed in voice pitch and  
216 the listeners were asked to indicate which *nine* was spoken by the person with the higher  
217 pitch. The first digit always had an  $f_0$  that matched one of the GPRs from the main experiment  
218 (i.e., 79, 150, 285 Hz). All GPRs from the main experiment were presented in separate  
219 staircases. The VTL was fixed at 12 cm in all three staircases. The second digit differed by delta  
220 cents from the first digit with an initial difference of 100 cents (i.e., one semitone). The  
221 direction of this difference was randomized in each trial. For the first four reversals in the  
222 direction of the staircase, the pitch difference was decreased by 10 cents following a correct  
223 response and increased by 30 cents following an incorrect response. From the fifth reversal  
224 onward, the step sizes were 2 and 6 cents for down- and up-steps, respectively. Each staircase  
225 was terminated after the twelfth reversal and the JND for GPR was defined as the arithmetic  
226 mean of delta cents visited on all reversal trials after the fifth reversal. Finally, JNDs were  
227 averaged across all three staircases.

228 For VTL-JNDs, we used a similar procedure. On each trial, two versions of the spoken  
229 digit nine were presented that differed in vocal timbre. The first digit always had a spectral  
230 envelope that matched one of the VTLs from the main experiment. All VTLs from the main  
231 experiment (i.e., 8, 12, 18 cm) were recycled in separate staircases, and the GPR was fixed at  
232 150 Hz in all three staircases. The difference in VTL was realized as spectral envelope ratio  
233 (SER), and each staircase started with an initial SER of 12 %. The SER was manipulated using  
234 up- and down-steps of 3 % and 1 % for the first four reversals, and 0.6 % and 0.2 % from the  
235 fifth reversal onward. Since VTL information has been associated with the perception of talker  
236 size (Smith et al., 2005), we asked listeners to indicate which *nine* was spoken by the smaller  
237 person (similar to Roswadowitz et al., 2014). The VTL-JND was defined as the arithmetic  
238 mean of SERs visited on all second-phase reversal trials.

## 239 **E. Data analysis**

240 Raw data were prepared for statistical analysis using MATLAB. We first calculated listeners'  
241 accuracies (i.e., proportion of correctly reported target digits) per condition, experiment, and  
242 digit position (i.e., the four digits per trial). To investigate the effect of continuity in different  
243 target voice features on cocktail-party listening, we calculated separate scores for the benefits  
244 from continuity in VTL & GPR, VTL-only continuity, and GPR-only continuity. To do this, we first  
245 logit-transformed accuracies separately for each listener, experiment, condition, and digit  
246 position. To correct for values of 0 and 1, we used an approach established in signal detection  
247 theory (Macmillan and Creelman, 2005), where  $p_{\text{correct}} = 1$  was set to  $p_{\text{correct}} = 1 - 1/(2N)$ , and  
248  $p_{\text{correct}} = 0$  was set to  $p_{\text{correct}} = 1/2N$  ( $N$  is the number of responses that entered the average; in  
249 this case,  $N=36$ ) (similar to Hartwigsen et al., 2015). Finally, we calculated the difference of the

250 logit-transformed accuracies in the GPR & VTL change versus each of the other three  
251 conditions. For example, the difference between GPR change and GPR & VTL change  
252 conditions quantified the benefit from continuity in VTL only, because VTL was fixed in the  
253 GPR change condition but varied in the GPR & VTL change condition (Fig. 1B). Similarly, we  
254 quantified the benefits from continuity in GPR only (VTL change – GPR & VTL change), and the  
255 benefits from continuity in VTL & GPR (Fixed voice – GPR & VTL change).

256 Previous work has shown that the previous-digit-correct benefit (PDCB) is a sensitive  
257 measure to capture benefits that arise from perceptual voice continuity (Bressler et al., 2014).  
258 The PDCB relates the probabilities of being correct on the current digit conditioned on  
259 whether or not the previous digit was correctly reported [ $P(C_{i-1})$  vs.  $P(NC_{i-1})$ ]. Like Bressler  
260 et al. (2014), we calculated the PDCB as the natural logarithm of the ratio of these conditional  
261 probabilities. For the calculation of both types of conditional probabilities, we used the same  
262 correction formula that we applied to the proportion-correct values. PDCBs were calculated  
263 separately for each listener, experiment, and condition.

264 A PDCB of zero would indicate that being correct on the previous digit had no effect  
265 on the probability of being correct on the current digit. The greater the PDCB, the greater the  
266 benefit from having been correct on the previous digit. If listeners were better at tracking the  
267 target stream with continuity in certain voice features, we should observe a greater PDCB in  
268 conditions in which voice features were kept constant in the target stream compared to  
269 conditions in which they varied across digits. Hence, investigation of the PDCB allowed us to  
270 explore whether the listeners' abilities to tune into the target stream are modulated by  
271 continuity in certain voice features.

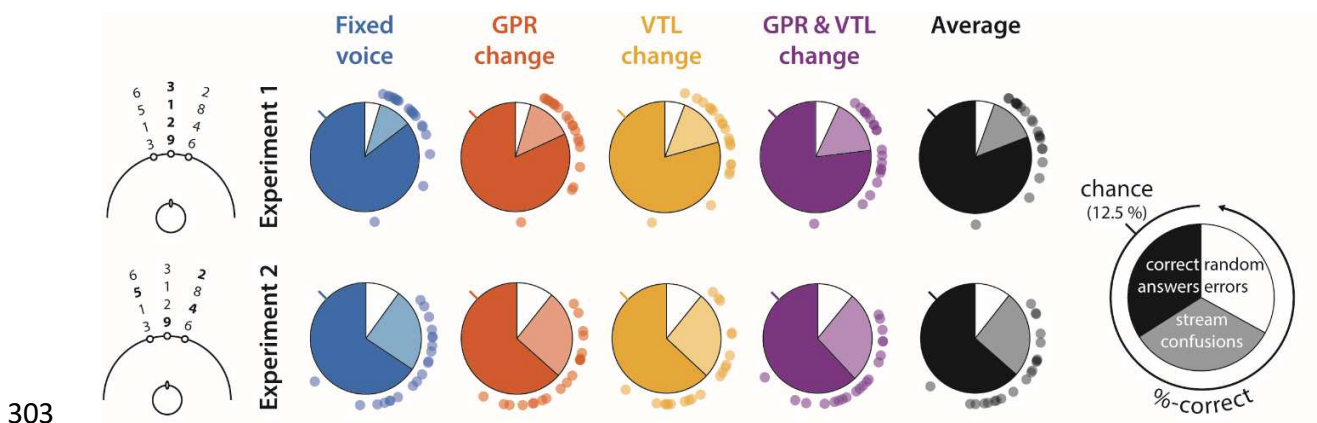
272           The statistical analyses were performed in R (R Core Team, 2017) using RStudio  
273 (version 1.1.383). Linear mixed-effects models as implemented in the *lme4* package (Bates et  
274 al., 2015) were fitted separately to continuity benefits and PDCBs. In all model fits, we  
275 followed an iterative procedure: starting with the intercept-only models, we first added fixed-  
276 and then random-effects terms in a stepwise fashion. After each step, we fitted the model  
277 using maximum-likelihood estimation, and assessed the change in model fit using likelihood-  
278 ratio tests.

279           To investigate the potential effects of Continuity type (VTL & GPR, VTL-only, GPR-only  
280 continuity) and Experiment (Exp. 1, Exp. 2) on continuity benefits, we modeled these  
281 predictors as fixed effects using deviation coding. To investigate the potential effect of Digit  
282 position (digit positions 1 to 4), we used backward difference coding; that is, we compared  
283 the continuity benefit for a given digit position to the benefit for the prior digit which allowed  
284 us to test for a successive increase in continuity benefits over digit positions. For the analysis  
285 of PDCBs, we investigated the potential effects of Condition (Fixed voice, GPR change, VTL  
286 change, GPR & VTL change) and Experiment using deviation coding. We derived p-values for  
287 individual model terms using the Satterthwaite approximation for degrees of freedom (Luke,  
288 2017). Post-hoc comparisons were performed using Tukey's range tests as implemented in  
289 the *lsmeans* package (Lenth, 2016). To provide an estimate of effect size for pairwise  
290 comparisons, we report unstandardized coefficients *b*.

### 291 III. Results

#### 292 A. Accuracy

293 Figure 2 shows the proportions of correctly reported target digits, stream confusions  
294 (reporting a digit from one of the two masker streams), and random errors (reporting a digit  
295 that was not present in the mixture) in each condition of both experiments as well as  
296 listeners' accuracies. All listeners performed well above chance (i.e., 0.125 or 1 out of 8  
297 possible response options) in all conditions of both experiments; when errors occurred,  
298 stream confusions were much more common than random errors in both Experiment 1 ( $t_{19} =$   
299  $9.13; p < 0.001; r = 0.90$ ) (Rosenthal and Rubin, 2003) and Experiment 2 ( $t_{19} = 9.56; p < 0.001; r =$   
300  $0.91$ ), even though there were more response options related to random errors (5) than  
301 stream confusions (2). Taken together, these results suggest that listeners were actively  
302 engaged in solving the cocktail-party problem.

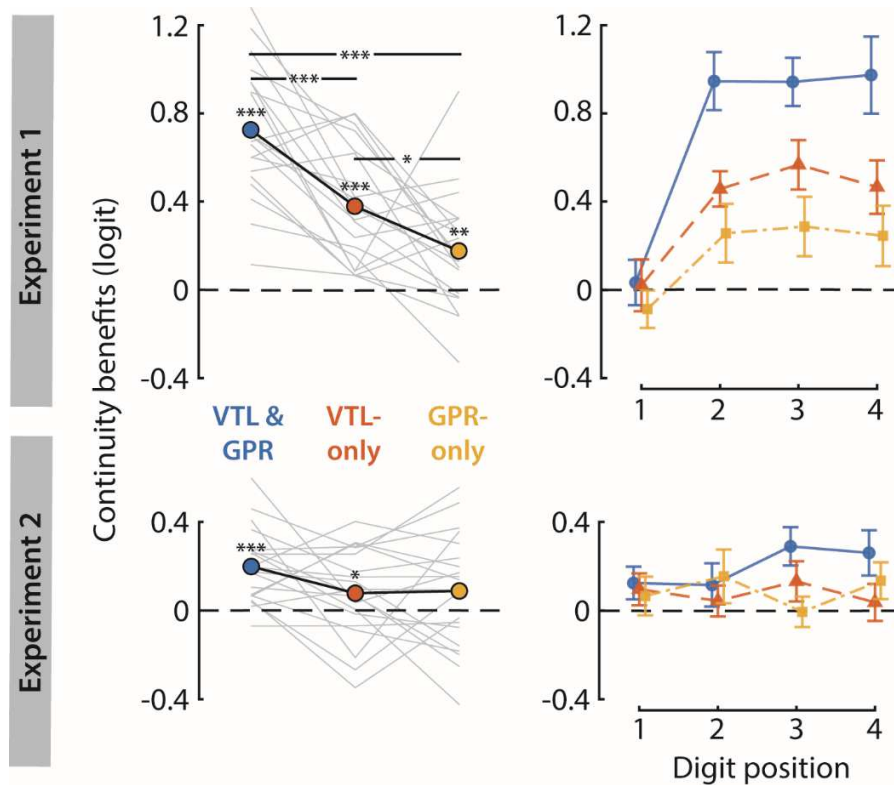


## 309 **B. Continuity benefits**

310 The main aim of the present study was to investigate the relative contributions of GPR and  
311 VTL to the voice continuity benefit. To quantify and compare the benefits from continuity in  
312 certain voice features, we calculated separate benefit scores for VTL & GPR, VTL-only as well as  
313 GPR-only continuity (see *Data analysis* for details).

314 One-sample *t*-tests revealed that listeners benefited significantly from all three  
315 continuity types (VTL & GPR:  $t_{159} = 9.58$ ;  $p < 0.001$ ;  $r = 0.61$ ; VTL-only:  $t_{159} = 6.19$ ;  $p < 0.001$ ;  $r =$   
316  $0.44$ ; GPR-only:  $t_{159} = 2.64$ ;  $p = 0.009$ ;  $r = 0.20$ ) when the target location was kept constant  
317 across digits (Exp. 1) (Fig. 3, top left). When the target location varied from digit to digit (Exp.  
318 2), listeners benefited significantly from continuity in VTL & GPR ( $t_{155} = 4.55$ ;  $p < 0.001$ ;  $r = 0.34$ ).  
319 The benefits from VTL-only and GPR-only continuity were similar in size, but only VTL-only  
320 continuity reached significance ( $t_{155} = 2.17$ ;  $p = 0.031$ ;  $r = 0.17$ ; GPR-only continuity:  $t_{155} = 1.84$ ;  
321  $p = 0.068$ ;  $r = 0.15$ ) (Fig. 3, bottom left).





322

323 **Figure 3.** Benefits from continuity in VTL & GPR, VTL-only, and GPR-only in Experiments 1 (top row) and  
 324 row). The left-hand side of the figure shows continuity benefits averaged across digit positions. Light gray lines  
 325 show continuity benefits for each individual listener, black lines show the mean across listeners. Significant  
 326 benefits are denoted by the asterisks directly above the colored dots. Significant differences across continuity  
 327 types are denoted by the asterisks within lines. \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ . The right-hand side of the  
 328 figure shows the continuity benefits as a function of digit position. Symbols show mean continuity benefits; error  
 329 bars denote the standard errors of the means.

330

331 In addition to these findings, several basic observations can be made by visual  
 332 inspection of Figure 3: first, spatial discontinuity in Experiment 2 greatly reduced overall  
 333 voice-feature continuity benefits (top vs. bottom row); second, the benefit scores decreased  
 334 considerably across the three continuity types in Experiment 1 (top left); third, the continuity  
 335 benefits emerged rapidly from the first to the second digit in Experiment 1 (top right).

336 These observations were confirmed by fitting linear mixed-effects models to the  
 337 benefit scores. The best-fitting model included the interaction terms between the factors  
 338 Continuity type and Experiment ( $F_{2,860.69} = 8.99$ ;  $p < 0.001$ ), and between the factors Digit

339 position and Experiment ( $F_{3,862.96} = 10.58; p < 0.001$ ) as well as all three main factors Continuity  
340 type ( $F_{2,860.69} = 21.42; p < 0.001$ ), Digit position ( $F_{3,20.92} = 6.71; p = 0.002$ ), and Experiment ( $F_{1,19.29} =$   
341 23.93;  $p < 0.001$ ) as fixed effects. The random-effects terms included the subject-specific  
342 random intercepts as well as the subject-specific random slopes for the factors Experiment  
343 and Digit Position.

344 To explore the Experiment-by-Continuity type interaction, we performed pairwise  
345 comparisons between all combinations of continuity types in each experiment. In Experiment  
346 1, all pairwise comparisons revealed significant differences across benefit scores (Fig. 3, top  
347 left): the listeners benefited more from continuity in VTL & GPR than from continuity in either  
348 VTL alone ( $t_{860.69} = 4.77; p < 0.001; b = 0.3462$ ) or GPR alone ( $t_{860.69} = 7.57; p < 0.001; b = 0.5489$ ).  
349 These results suggest that the effects of VTL and GPR continuity were additive and that  
350 listeners exploited all of the continuity available in the target stream instead of focusing on a  
351 single voice feature. Importantly, however, the results showed greater benefits from VTL-only  
352 compared to GPR-only continuity ( $t_{860.69} = 2.79; p = 0.015; b = 0.2027$ ), suggesting that  
353 perceptual grouping of target digits relied more on continuity in VTL than GPR. In Experiment  
354 2, none of the pairwise comparisons between continuity types turned out to be significant ( $p$   
355  $\geq 0.282$ ) (Fig. 3, bottom left).

356 Next, we explored the Experiment-by-Digit position interaction by performing  
357 pairwise comparisons between all combinations of digit positions in each experiment. In  
358 Experiment 1, we found significant differences in continuity benefits for the comparisons  
359 between digit position 1 and all other digit positions (Digit 1 vs. Digit 2:  $t_{53.60} = -5.66; p < 0.001;$   
360  $b = -0.5648$ ; Digit 1 vs. Digit 3:  $t_{42.51} = -5.95; p < 0.001; b = -0.6104$ ; Digit 1 vs. Digit 4:  $t_{32.25} = -$

361 4.62;  $p < 0.001$ ;  $b = -0.5731$ ). None of the other pairwise comparisons yielded significant  
 362 differences ( $p \geq 0.968$ ). In Experiment 2, continuity benefits did not differ significantly for any  
 363 pairwise comparison between digit positions ( $p \geq 0.962$ ). Taken together, these results  
 364 showed a rapid emergence of continuity benefits (from the first to the second digit) without  
 365 any further significant increase at later digit positions. This rapid emergence of continuity  
 366 benefits was evident in Experiment 1 (Fig. 3, top right) but not in Experiment 2 (Fig. 3, bottom  
 367 right), suggesting that it depended on spatial continuity. The results for the effects of the  
 368 three main factors Continuity type, Experiment, and Digit position are summarized in Table I.

369  
 370 **Table I.** Continuity benefits: results for the effects of Continuity type, Experiment, and Digit  
 371 position.  
 372

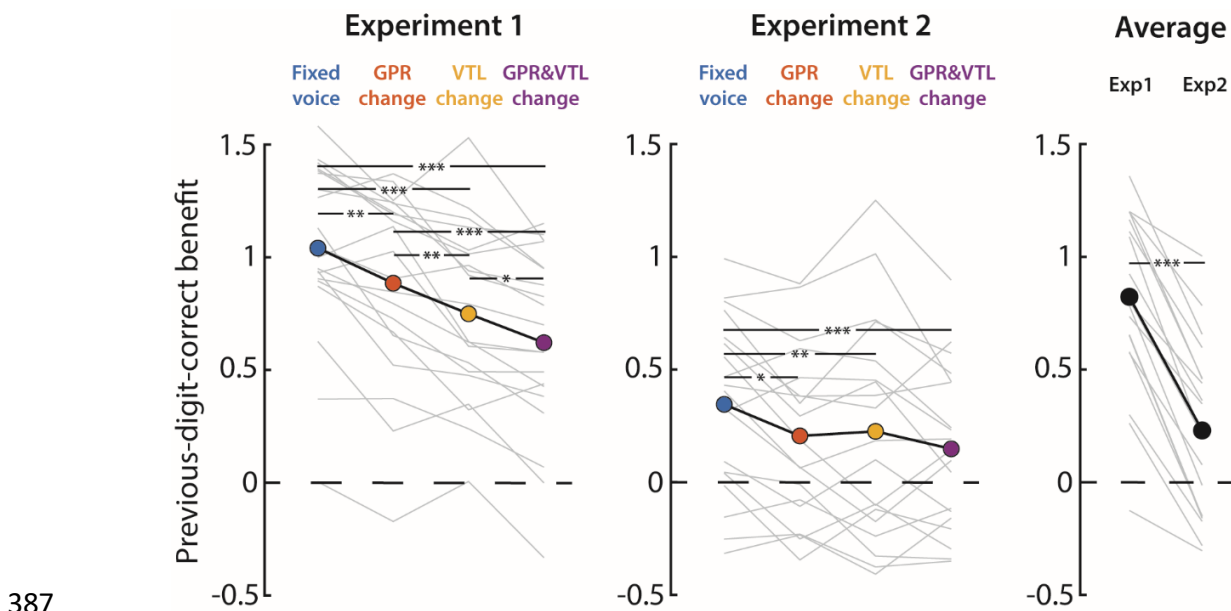
Effect	b	t	df	p
<b>Continuity type</b>				
VTL & GPR continuity vs. grand mean	0.3715	6.23	860.70	<b>&lt; 0.001</b>
VTL-only continuity vs. grand mean	-0.0823	-1.38	860.70	0.168
<b>Experiment</b>				
Experiment 1 vs. Experiment 2	0.3055	4.89	19.30	<b>&lt; 0.001</b>
<b>Digit position</b>				
Digit 2 vs. Digit 1	0.4316	3.57	23.10	<b>0.002</b>
Digit 3 vs. Digit 2	0.0852	0.53	21.70	0.600
Digit 4 vs. Digit 3	-0.0339	-0.22	19.20	0.829

373 P-values for significant comparisons are marked by bold face.  
 374

### 375 **C. Previous-digit-correct benefit**

376 To investigate whether listeners' abilities to tune into the target stream are modulated by  
 377 voice-feature continuity, we calculated the PDCB (similar to Bressler et al., 2014) which relates  
 378 the probability of being correct on the current digit conditioned on whether or not the

379 listener was correct on previous digit. One-sample  $t$ -tests revealed significant PDCBs in all  
 380 conditions of both experiments ( $p < 0.05$ ). Furthermore, the data shown in Figure 4 suggest  
 381 that the PDCBs differed across conditions (especially in Experiment 1) and that listeners  
 382 attained greater overall PDCBs in Experiment 1 than Experiment 2. Indeed, the best-fitting  
 383 model included the interaction term between the factors Experiment and Condition ( $F_{3,11330} =$   
 384  $6.99; p < 0.001$ ) as well as the main factors Experiment ( $F_{1,19} = 145.71; p < 0.001$ ) and Condition  
 385 ( $F_{3,11330} = 37.53; p < 0.001$ ) as fixed effects. The random effects were the subject-specific  
 386 random intercepts and the subject-specific random slopes for the factor Experiment.



388 **Figure 4.** Previous-digit-correct benefits (PDCBs) shown for each condition in Experiments 1 (left) and 2 (middle)  
 389 as well as averaged across conditions in each experiment (right). Light gray lines show performance for each  
 390 individual listener, black lines show the mean across listeners. Significant differences across conditions and  
 391 experiments are denoted by asterisks. \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

392  
 393 We explored the Experiment-by-Condition interaction by performing pairwise  
 394 comparisons for all combinations of conditions in each experiment. In both experiments, the  
 395 PDCBs were greater in the *Fixed voice* condition compared to all other conditions (Fig. 4,  
 396 'Experiment 1' and 'Experiment 2'; see Table II for details), showing that the listeners were less

397 able to tune into the target stream when the target voices changed in either GPR, VTL, or  
 398 both. Importantly, in Experiment 1 but not in Experiment 2, we found significantly greater  
 399 PDCBs when the target voices changed in GPR compared to VTL (Fig. 4, left), showing that VTL  
 400 changes had a more detrimental effect on the ability to tune into the target stream in  
 401 Experiment 1.

402  
 403 **Table II.** Previous-digit-correct benefit: results of post-hoc comparisons for the Experiment-  
 404 by-Condition interaction. Degrees of freedom for all columns df = 11,329.89.  
 405

<b>Comparison</b>	<b>b</b>	<b>t</b>	<b>p</b>
<b>Experiment 1</b>			
Fixed voice vs. GPR change	0.1555	3.70	<b>0.001</b>
Fixed voice vs. VTL change	0.2912	6.92	<b>&lt; 0.001</b>
Fixed voice vs. GPR & VTL change	0.4193	9.96	<b>&lt; 0.001</b>
GPR change vs. VTL change	0.1357	3.22	<b>0.007</b>
GPR change vs. GPR & VTL change	0.2638	6.27	<b>&lt; 0.001</b>
VTL change vs. GPR & VTL change	0.1281	3.04	<b>0.013</b>
<b>Experiment 2</b>			
Fixed voice vs. GPR change	0.1531	3.59	<b>0.019</b>
Fixed voice vs. VTL change	0.1216	2.85	<b>0.023</b>
Fixed voice vs. GPR & VTL change	0.2040	4.78	<b>&lt; 0.001</b>
GPR change vs. VTL change	-0.0316	-0.74	0.881
GPR change vs. GPR & VTL change	0.0508	1.19	0.632
VTL change vs. GPR & VTL change	0.0824	1.93	0.214

406 P-values for significant comparisons are marked by bold face.  
 407

#### 408 **D. Just-noticeable differences**

409 The GPR and VTL values used in the present study were chosen to induce the perception of  
410 changes in talker identity. They correspond to a minimal difference of 1078 cents and 50 %  
411 spectral envelope ratio (SER), respectively. In the literature, about half of these step sizes have  
412 been reported to be sufficient to elicit the perception of a talker identity change (Gaudrain et  
413 al., 2009). Nevertheless, we checked whether all listeners were sensitive to the GPR and VTL  
414 changes in the two main experiments by comparing the minimal differences of our voice-  
415 feature manipulations to listeners' JNDs for GPR and VTL. The JNDs for GPR ranged from 12.33  
416 to 87.92 cents and were on average ( $M = 41.04$  cents) significantly smaller than the minimal  
417 GPR difference in the two main experiments (GPR:  $t_{19} = -213.33$ ;  $p < 0.001$ ;  $r = 1$ ). The same was  
418 also true for VTL-JNDs ( $M = 4.93$  % SER; ranging from 1.33 to 17.21 % SER) ( $t_{19} = -69.49$ ;  $p <$   
419  $0.001$ ;  $r = 1$ ).

420 Note that, expressed in average JNDs, the minimal difference between virtual talkers  
421 in the present study was larger for GPR (1078 cents correspond to about 26 JNDs) than VTL  
422 differences (50 % SER corresponds to about 10 JNDs). The perceptually larger change in GPR  
423 than VTL can therefore not explain our main finding that listeners benefited more from VTL  
424 than GPR continuity. Furthermore, individual listeners' JNDs for GPR and VTL were not  
425 correlated with listeners' benefits from continuity in the respective voice features (GPR:  $r_s =$   
426  $0.05$ ;  $p = 0.836$ ; VTL:  $r_s = 0.31$ ;  $p = 0.186$ ).

## 427 **IV. Discussion**

428 The present study investigated the effects of (dis-)continuity in two of the most salient voice  
429 features, GPR and VTL, on listeners' abilities to solve the cocktail-party problem. When the  
430 target location was fixed within a trial (Exp. 1), listeners showed the greatest benefits from  
431 continuity in both voice feature. The most important result, however, was that listeners  
432 showed greater benefits from continuity in VTL alone than GPR alone. Our results thus  
433 suggest that listeners used all the continuity available in the target stream, but when  
434 continuity was only available in one of the two voice features, VTL continuity was more  
435 beneficial for perceptual grouping.

### 436 **A. Different roles of VTL and GPR for grouping and segregation**

437 Our results might appear unexpected when juxtaposed to previous studies on the  
438 involvement of GPR and VTL in stream segregation (Darwin et al., 2003; Vestergaard et al.,  
439 2009). Notably, these studies manipulated the dissimilarity of target and masker streams in  
440 GPR and VTL and found that less dissimilarity in GPR than VTL was needed to yield  
441 comparable performance, suggesting that GPR is the more beneficial feature. A possible  
442 explanation for this apparent discrepancy is that the different experimental manipulations  
443 tapped into different aspects of cocktail-party listening: while manipulating the dissimilarity  
444 of competing streams in previous studies focused on the influence of GPR and VTL on  
445 *segregation*, the manipulation of voice-feature continuity for target speech in the present  
446 study allowed us to investigate the influence of GPR and VTL on *grouping*.

447 Theoretically, both segregation and grouping are important processes for cocktail-  
448 party listening as they lend support to the formation and selection of perceptual objects in

449 the auditory scene (for a recent review, see Shinn-Cunningham et al., 2017). However, GPR  
450 and VTL might contribute differently to these processes. For segregation, the listeners'  
451 differential sensitivity to GPR and VTL changes might play an important role. Consistent with  
452 previous studies (Ives et al., 2005; Smith et al., 2005), we found that a change in VTL had to be  
453 about twice as large as a change in GPR to be perceived by the listeners (4.93 % vs. 2.34 %). It  
454 is therefore not surprising that listeners are better at segregating two competing streams  
455 based on GPR compared to VTL differences, especially when these differences are small and  
456 no other perceptual features are available.

457 For grouping, however, listeners might rely on their experience with natural talkers. A  
458 natural talker's VTL is relatively fixed with only slight variations due to articulatory  
459 movements, whereas GPR varies considerably due to the use of prosodic cues in natural  
460 speech (Kania et al., 2006). Consequently, vocal tract features have been found to be more  
461 important for the identification of natural talkers than glottal fold features (Lavner et al.,  
462 2000). It is thus likely that the listeners in the present study benefited more from continuity in  
463 VTL because they have learned that VTL is the more reliable cue for the identification of  
464 natural talkers.

465 A potential caveat is that we only observed greater benefits from VTL than GPR  
466 continuity because of the specific values of GPR and VTL that were chosen. When GPR  
467 changed between consecutive target digits, the difference was at least 90 %. For VTL changes,  
468 we used a minimal difference of 50 %. These differences were chosen to elicit the perception  
469 of talker identity changes rather than variations within a talker's voice and were consistent  
470 with previous work showing that listeners perceive different talker identities at about half of



471 these magnitudes (Gaudrain et al., 2009). We did not assess whether the changes were indeed  
472 large enough to be perceived as separate talker identities with our specific stimuli, but we did  
473 confirm that all listeners were sensitive to these changes. Furthermore, the changes in GPR  
474 were perceptually (i.e., expressed in JNDs) larger than the changes in VTL. Also, individual  
475 sensitivity to GPR and VTL was not related to how much listeners benefited from continuity in  
476 the respective voice feature (i.e., there were no correlations between JNDs and voice-  
477 continuity benefits). It is thus unlikely that the specific GPR and VTL values used here can  
478 explain the greater benefits from VTL than GPR continuity.

479 Further support for a genuinely stronger contribution of VTL than GPR to perceptual  
480 grouping comes from a study on the phonemic restoration effect (Clarke et al., 2014). While  
481 phonemic restoration persisted changes in either voice feature, global speech intelligibility  
482 suffered more from VTL than GPR changes. Importantly, the GPR and VTL changes were  
483 comparable to the changes in the present study and listeners perceived them as a change in  
484 talker identity.

## 485 **B. Costs of spatial discontinuity**

486 A second aim of the present study was to investigate the effect of spatial discontinuity on  
487 listeners' abilities to group sounds based on voice-feature continuity. Introducing spatial  
488 discontinuity drastically reduced the benefits from voice-feature continuity. This finding can  
489 be explained in terms of a lack of temporal coherence across acoustic features (Shamma et al.,  
490 2011). When attention is allocated to a particular location, all other temporally coherent  
491 features (e.g., pitch and timbre) of the source at this location can be perceptually grouped. In  
492 Experiment 2, we broke the temporal coherence between location and voice features:

493 listeners had to divide spatial attention since they had no advance knowledge about the  
494 target location, which can explain why they benefited less from voice-feature continuity in  
495 Experiment 2.

496         The costs associated with spatial discontinuity were also evident in the evolution of  
497 continuity benefits over time. Listeners showed a large increase in continuity benefits from  
498 the first to the second target digit when they could maintain attention on one location.  
499 However, this rapid emergence of continuity benefits was lost when listeners had to switch  
500 spatial attention from one target digit to the next.

### 501 **C. Voice-feature continuity as a perceptually driven bias of selective** 502 **attention**

503 These results shed light on the temporal dynamics of selective attention and support the  
504 notion that attention operates on perceptual objects (Shinn-Cunningham, 2008). Obviously,  
505 there were no continuity benefits for the first digit within a trial, but as long as listeners could  
506 maintain selective attention on the same location, they latched onto whatever voice feature  
507 was continuous across subsequent target digits. Such a rapid emergence of continuity  
508 benefits is remarkable given that the build-up of selective attention can take up to a couple of  
509 seconds (Cusack et al., 2004). It is difficult to imagine that listeners volitionally decided to  
510 focus their attention on a specific voice feature, particularly because they did not know in  
511 advance which, if any, voice feature would be continuous in the target stream.

512         Our results can be rather interpreted in terms of a perceptually driven bias of selective  
513 attention (Bressler et al., 2014): once listeners had encoded certain voice features from the  
514 first digit, continuity in any of these features might have biased the listeners to focus on these

515 features in subsequent digits. This explanation does not rely on a rather slow build-up of  
516 selective attention; instead, it is based on the assumption that whatever voice feature is in the  
517 attentional foreground of the current digit will be perceptually enhanced in the mixture of  
518 subsequent digits.

519         If the above conjecture is true, then listeners should have only benefited from  
520 continuity in a certain voice feature once this feature was already in their attentional focus. In  
521 other words, listeners should have been more likely to correctly report the current target digit  
522 if they had correctly reported the previous target digit and this benefit should be greater  
523 when voice features were continuous across target digits. Our results on the PDCB showed  
524 that this was indeed the case: the benefits from being correct on the previous digit were  
525 greater when both GPR and VTL were continuous in the target stream compared to when  
526 either one or both voice features changed, showing that continuity in both GPR and VTL  
527 helped listeners to direct attention to the next target digits. This finding was independent  
528 from spatial (dis-)continuity; however, when listeners knew where the next target would  
529 appear, they were generally better at tuning into the target stream. Furthermore, with spatial  
530 continuity, listeners were more likely to tune into the target stream based on VTL than GPR  
531 continuity. Together with the greater continuity benefit from VTL than GPR continuity, this  
532 result provides converging evidence for the importance of VTL for perceptual grouping.

#### 533 **D. Implications for cochlear-implant users**

534 Our results are not only informative about the use of GPR and VTL for perceptual grouping in  
535 normal-hearing listeners, they also have implications for cocktail-party listening in cochlear-  
536 implant (CI) users. Cocktail-party listening is severely impaired in CI users (e.g., Loizou et al.,

537 2009). This is likely due to the reduced spectral resolution of the implant which hinders the  
538 analysis of voice features (Stickney et al., 2004). Specifically, it has been shown that CI users  
539 benefit much less from talker differences between target and masker speech than normal-  
540 hearing listeners and that this is even the case when target and masker talkers differ in sex.  
541 Furthermore, while normal-hearing listeners make use of both GPR and VTL differences for  
542 talker sex categorization, CI users seem to rely exclusively on differences in GPR (Fuller et al.,  
543 2014; Meister et al., 2016) which has been attributed to their limited access to VTL cues  
544 (Gaudrain and Başkent, 2018).

545 It remains an open question to what extent, if at all, CI users can benefit from  
546 continuity in a single voice feature in the cocktail party. However, the relative importance of  
547 VTL continuity for perceptual grouping found in the present study together with previous  
548 findings suggest that CI users often fail to solve the cocktail-party problem because of  
549 impaired processing of VTL information.

## 550 **E. A potential neural mechanism for dealing with voice-feature changes in** 551 **the cocktail party**

552 Relatively little is known about the neural mechanisms supporting perceptual grouping in the  
553 cocktail party. Yet, there is evidence that changes in talker sex (Shomstein and Yantis, 2006)  
554 and pitch (Hill and Miller, 2009; Lee et al., 2013) of a target sound are processed in bilateral  
555 areas of temporal cortices. Furthermore, activation in parts of these areas (i.e., left mid-  
556 posterior superior temporal gyrus) has been found to correlate with the comprehension of  
557 speech in noise (Evans et al., 2016), suggesting that it is behaviorally relevant for cocktail-  
558 party listening.

559           Neuroimaging work using clear speech suggests that robust speech comprehension in  
560 the context of GPR and VTL changes relies on functional interactions between left- and right-  
561 hemispheric areas that are sensitive to glottal-fold and vocal-tract information (Kreitewolf et  
562 al., 2014; von Kriegstein et al., 2010). These are areas in left and right Heschl's gyri that process  
563 glottal fold information relevant for the recognition of linguistic prosody and vocal pitch, as  
564 well as left and right posterior superior temporal areas that process vocal tract information  
565 relevant for the recognition of phonemes and vocal timbre. It is possible that these functional  
566 interactions are also at play when dealing with GPR and VTL changes in the cocktail party.

## 567 **V. Conclusion**

568 The present findings show that continuity in voice features helps perceptual grouping  
569 potentially because target voice features guide selective attention in the cocktail party. Most  
570 importantly, however, we found that listeners' abilities to solve the cocktail-party problem  
571 benefit more from continuity in VTL than GPR. This is likely a result of the differential  
572 importance of VTL and GPR for the identification of natural talkers: listeners might rely more  
573 on VTL continuity for perceptual grouping because they have learned that a natural talker's  
574 VTL is effectively fixed. Furthermore, these results might explain why cochlear-implant users,  
575 who have reduced access to VTL cues, particularly struggle in the cocktail party.

## 576 **Acknowledgments**

577 This work was supported by an Erasmus Mundus exchange stipend to JK and a Quebec  
578 Research Scholar grant to MS. We thank Scott Bressler for providing us with the original  
579 spoken digits from a previous study (Bressler et al., 2014), Sarah Tune for her help in setting  
580 up the mixed-effects models in R, and Malte Wöstmann for comments on an earlier version of  
581 this manuscript.

## 582 **References**

583 Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). "Fitting Linear Mixed-Effects  
584 Models Using lme4," *Journal of Statistical Software* **67**(1), pp. 1–48.

585 Baumann, O., and Belin, P. (2010). "Perceptual scaling of voice identity: common  
586 dimensions for different vowels and speakers," *Psychological Research* **74**(1), pp. 110–120.

587 Best, V., Ozmeral, E. J., Kopčo, N., and Shinn-Cunningham, B. G. (2008). "Object  
588 continuity enhances selective auditory attention," *Proceedings of the National Academy of  
589 Science USA* **105**(35), pp. 13174–13178.

590 Bregman, A. S. (1990). *Auditory scene analysis: The perceptual organization of sound* (MIT  
591 press, Cambridge, Massachusetts)

592 Bressler, S., Masud, S., Bharadwaj, H., and Shinn-Cunningham, B. (2014). "Bottom-up  
593 influences of voice continuity in focusing selective auditory attention," *Psychological research*  
594 **78**(3), pp. 349–360.

595 Brungart, D. S., and Simpson, B. D. (2007). "Cocktail party listening in a dynamic  
596 multitalker environment," *Perception & Psychophysics* **69**(1), pp. 79–91.

597 Cherry, E. C. (1953). "Some experiments on the recognition of speech, with one and with  
598 two ears," *The Journal of the Acoustical Society of America* **25**(5), pp. 975–979.

599 Clarke, J., Gaudrain, E., Chatterjee, M., and Başkent, D. (2014). "T'ain't the way you say it,  
600 it's what you say—Perceptual continuity of voice and top–down restoration of speech," *Hearing  
601 research* **315**, pp. 80–87.

602 Cusack, R., Decks, J., Aikman, G., and Carlyon, R. P. (2004). "Effects of location, frequency  
603 region, and time course of selective attention on auditory scene analysis," *Journal of*  
604 *experimental psychology: human perception and performance* **30**(4), pp. 643–656.

605 Darwin, C. J., Brungart, D. S., and Simpson, B. D. (2003). "Effects of fundamental  
606 frequency and vocal-tract length changes on attention to one of two simultaneous talkers," *The*  
607 *Journal of the Acoustical Society of America* **114**(5), pp. 2913–2922.

608 Evans, S., McGettigan, C., Agnew, Z. K., Rosen, S., and Scott, S. K. (2016). "Getting the  
609 cocktail party started: masking effects in speech perception," *Journal of cognitive neuroscience*  
610 **28**(3), pp. 483–500.

611 Fuller, C. D., Gaudrain, E., Clarke, J. N., Galvin, J. J., Fu, Q. J., Free, R. H., and Başkent, D.  
612 (2014). "Gender categorization is abnormal in cochlear implant users," *Journal of the*  
613 *Association for Research in Otolaryngology* **15**(6), pp. 1037–1048.

614 Gaudrain, E., Li, S., Ban, V. S., and Patterson, R. D. (2009). "The role of glottal pulse rate  
615 and vocal tract length in the perception of speaker identity," *Interspeech 2009*, Brighton, pp.  
616 148–151.

617 Gaudrain, E., and Başkent, D. (2018). "Discrimination of voice pitch and vocal-tract  
618 length in cochlear implant users," *Ear and hearing* **39**(2), pp. 226–237.

619 Hartwigsen, G., Golombek, T., and Obleser, J. (2015). "Repetitive transcranial magnetic  
620 stimulation over left angular gyrus modulates the predictability gain in degraded speech  
621 comprehension," *Cortex* **68**, pp. 100–110.

622 Hill, K. T., and Miller, L. M. (2009). "Auditory attentional control and selection during  
623 cocktail party listening," *Cerebral cortex* **20**(3), pp. 583–590.



624 Ives, D. T., Smith, D. R., and Patterson, R. D. (2005). "Discrimination of speaker size from  
625 syllable phrases," *The Journal of the Acoustical Society of America* **118**(6), pp. 3816–3822.

626 Kaernbach, C. (1991). "Simple adaptive testing with the weighted up-down method,"  
627 *Attention, Perception, & Psychophysics* **49**(3), pp. 227–229.

628 Kania, R. E., Hartl, D. M., Hans, S., Maeda, S., Vaissiere, J., and Brasnu, D. F. (2006).  
629 "Fundamental frequency histograms measured by electroglottography during speech: a pilot  
630 study for standardization," *Journal of voice* **20**(1), pp. 18–24.

631 Kawahara, H., Morise, M., Takahashi, T., Nisimura, R., Irino, T., and Banno, H. (2008).  
632 "TANDEM-STRAIGHT: A temporally stable power spectral representation for periodic signals  
633 and applications to interference-free spectrum, F0, and aperiodicity estimation," *ICASSP 2008*,  
634 Las Vegas, pp. 3933–3936.

635 Kidd Jr, G., Arbogast, T. L., Mason, C. R., and Gallun, F. J. (2005). "The advantage of  
636 knowing where to listen," *The Journal of the Acoustical Society of America* **118**(6), pp. 3804–  
637 3815.

638 Kitterick, P. T., Bailey, P. J., and Summerfield, A. Q. (2010). "Benefits of knowing who,  
639 where, and when in multi-talker listening," *The Journal of the Acoustical Society of America*  
640 **127**(4), pp. 2498–2508.

641 Kreitewolf, J., Gaudrain, E., and von Kriegstein, K. (2014). "A neural mechanism for  
642 recognizing speech spoken by different speakers," *Neuroimage* **91**, pp. 375–385.

643 Larson, E., and Lee, A. K. (2013). "Influence of preparation time and pitch separation in  
644 switching of auditory attention between streams," *The Journal of the Acoustical Society of*  
645 *America* **134**(2), EL165-EL171.

646 Lee, A. K., Rajaram, S., Xia, J., Bharadwaj, H., Larson, E., Hämäläinen, M., and Shinn-  
647 Cunningham, B. G. (2013). "Auditory selective attention reveals preparatory activity in  
648 different cortical regions for selection based on source location and source pitch," *Frontiers in*  
649 *neuroscience* **6**, 190.

650 Lavner, Y., Gath, I., and Rosenhouse, J. (2000). "The effects of acoustic modifications on  
651 the identification of familiar voices speaking isolated vowels," *Speech Communication* **30**(1),  
652 pp. 9–26.

653 Lenth, R. V. (2016). "Least-squares means: the R package lsmeans," *Journal of*  
654 *Statistical Software* **69**(1), pp. 1–33.

655 Loizou, P. C., Hu, Y., Litovsky, R., Yu, G., Peters, R., Lake, J., and Roland, P. (2009).  
656 "Speech recognition by bilateral cochlear implant users in a cocktail-party setting," *The*  
657 *Journal of the Acoustical Society of America* **125**(1), pp. 372–383.

658 Luke, S. G. (2017). "Evaluating significance in linear mixed-effects models in R,"  
659 *Behavior Research Methods* **49**, pp. 1494–1502.

660 Macmillan, N. A., and Creelman, C. D. (2005). *Detection Theory: A User's Guide (2nd Ed)*  
661 (Cambridge University Press, Cambridge, UK)

662 Mathias, S. R., and von Kriegstein, K. (2014). "How do we recognise who is speaking,"  
663 *Front Biosci (Schol Ed)* **6**, pp. 92–109.

664 Meister, H., Fürsen, K., Streicher, B., Lang-Roth, R., and Walger, M. (2016). "The use of  
665 voice cues for speaker gender recognition in cochlear implant recipients," *Journal of Speech,*  
666 *Language, and Hearing Research* **59**(3), pp. 546–556.

667 R Core Team (2017). *R: A Language and Environment for Statistical Computing* (Vienna:  
668 The R Foundation for Statistical Computing)

669 Rosenthal, R., and Rubin, D. B. (2003). "r equivalent: A simple effect size indicator,"  
670 *Psychological methods* **8**(4), pp. 492–496.

671 Roswadowitz, C., Mathias, S. R., Hintz, F., Kreitewolf, J., Schelinski, S., and von  
672 Kriegstein, K. (2014). „Two cases of selective developmental voice-recognition impairments,"  
673 *Current Biology* **24**(19), pp. 2348–2353.

674 Shamma, S. A., Elhilali, M., and Micheyl, C. (2011). „Temporal coherence and attention  
675 in auditory scene analysis," *Trends in neurosciences* **34**(3), pp. 114–123.

676 Shinn-Cunningham, B. G. (2008). "Object-based auditory and visual attention," *Trends*  
677 *in Cognitive Science* **12**(5), pp. 182–186.

678 Shinn-Cunningham, B. G., Best, V., and Lee, A. K. (2017). *Auditory Object Formation and*  
679 *Selection. In: The Auditory System at the Cocktail Party* (Springer, Cham), pp. 7–40.

680 Shomstein, S., and Yantis, S. (2006). "Parietal cortex mediates voluntary control of  
681 spatial and nonspatial auditory attention," *Journal of Neuroscience* **26**(2), pp. 435–439.

682 Smith, D. R., Patterson, R. D., Turner, R., Kawahara, H., and Irino, T. (2005). "The  
683 processing and perception of size information in speech sounds," *The Journal of the*  
684 *Acoustical Society of America* **117**(1), pp. 305–318.

685 Stickney, G. S., Zeng, F. G., Litovsky, R., and Assmann, P. (2004). "Cochlear implant  
686 speech recognition with speech maskers," *The Journal of the Acoustical Society of America*  
687 **116**(2), pp. 1081–1091.

688 Vestergaard, M. D., Fyson, N. R., and Patterson, R. D. (2009). "The interaction of vocal  
689 characteristics and audibility in the recognition of concurrent syllables," *The Journal of the*  
690 *Acoustical Society of America* **125**(2), pp. 1114–1124.

691 von Kriegstein, K., Smith, D. R., Patterson, R. D., Kiebel, S. J., and Griffiths, T. D. (2010).  
692 "How the human brain recognizes speech in the context of changing speakers," *Journal of*  
693 *Neuroscience* **30**(2), pp. 629–638.

694 Zwicker, I. E., and Fastl, I. H. (1999). *Loudness. In: Psychoacoustics* (Springer, Berlin  
695 Heidelberg).