

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/94945/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Culling, John F. 2016. Speech intelligibility in virtual restaurants. *Journal of the Acoustical Society of America* 140 (4) , 2418. 10.1121/1.4964401

Publishers page: <http://dx.doi.org/10.1121/1.4964401>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



# Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <http://orca.cf.ac.uk/94945/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Culling, John Francis 2016. Speech intelligibility in virtual restaurants. Journal of the Acoustical Society of ASmerica Item availability restricted. file

Publishers page:

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1 **Speech intelligibility in virtual restaurants.**

2 *John F. Culling, School of Psychology, Cardiff University, Tower Building, Park*  
3 *Place, Cardiff, CF10 3AT, U.K.*

4

5 Date: Monday, 12 September 2016

6 Running title: Virtual restaurant.

7 PACS numbers:

8

9 Correspondence address: -

10 Prof John Culling  
11 School of Psychology,  
12 Cardiff University,  
13 Tower Building, Park Place,  
14 Cardiff,  
15 CF10 3AT  
16 U.K.  
17 CullingJ@cf.ac.uk

18

19

20

**Abstract.**

21 Speech reception thresholds (SRTs) for a target voice on the same virtual table were  
22 measured in various restaurant simulations under conditions of masking by between 1  
23 and 8 interferers at other tables. Results for different levels of reverberation and  
24 different simulation techniques were qualitatively similar. SRTs increased steeply  
25 with the number of interferers, reflecting progressive failure to perceptually unmask  
26 the target speech as the acoustic scene became more complex. For a single interferer,  
27 continuous noise was the most effective masker, and a single interfering voice of  
28 either gender was least effective. With two interferers, evidence of informational  
29 masking emerged as a difference in SRT between forward and reversed speech, but  
30 SRTs for all interferer types progressively converged at 4 and 8 interferers. In  
31 simulation based on a real room, this occurred at a signal-to-noise ratio of around -5  
32 dB.

## 33 I. INTRODUCTION

34           Speech intelligibility in noise has been studied intensively in the laboratory  
35 using stimuli that varied widely in their ecological validity, but few have attempted to  
36 fully recreate a realistic listening experience. Early studies were limited by the  
37 technology of the day and generally presented words, non-words or sentence materials  
38 against white noise or pure tones (Miller, 1947; Licklider, 1948), high-/low-pass  
39 filtered noise (Fletcher and Galt, 1950) or modulated noise (Miller, 1947). These  
40 studies provided insights into the way that basic mechanisms of masking and hearing  
41 can contribute to the understanding of speech. More recent experiments have  
42 introduced realistic binaural cues (Bronkhorst and Plomp, 1988), multiple interfering  
43 sources (Hawley et al., 2004), room reverberation (Beutelmann and Brand, 2006) and  
44 the combination of all three (Culling, 2013; Westermann and Buchholz, 2015). The  
45 importance of these developments is that realistic, but experimentally controlled  
46 stimuli enable us to determine the roles of different mechanisms in real life. The  
47 present experiment addressed two questions in particular. The relative roles of  
48 informational and energetic masking and the speech-to-noise ratios (SNRs) that can  
49 occur in real-world listening.

50           Informational masking has been a topic of intense interest over the last 15  
51 years. Under some circumstances, listeners can fail to understand speech in conditions  
52 where conventional (“energetic”) masking mechanisms would be expected to have  
53 little role. For instance, Brungart et al. (2001) found that the intelligibility of  
54 sentences containing color/number combinations could be substantially lower when  
55 masked by similar sentences than when masked by noise whose spectral content and  
56 modulation were matched to the masking sentences. The lower intelligibility was  
57 attributed to the addition of informational masking. On one hand, the listening

58 situation was very unrealistic, in that the sentences were highly stylized and  
59 interfering sentences were saying very similar things to the target sentences. On the  
60 other hand, it can be argued that the traditional use of noise is unrealistic and that  
61 interfering speech is a more typical form of masking in everyday life. The question  
62 therefore arises, of whether informational masking has a prominent role in those  
63 everyday life situations where listening becomes difficult.

64         The second question concerns what those difficult everyday life situations  
65 would be. In laboratory studies, speech reception thresholds for 50% intelligibility  
66 (SRTs) can be extremely low under some circumstances. When interfering noise is  
67 strongly modulated SRTs can reach -23 dB in speech-shaped noise (Rhebergen and  
68 Versfeld, 2005). When spatial configurations are favorable, SRTs of around -12 dB  
69 have been reported for a continuous speech-shaped noise interferer and -20 dB for a  
70 speech interferer (Hawley et al. 2004). This advantage for a speech interferer is partly  
71 attributable to the modulation of the speech, but probably also to the harmonic  
72 structure of its voiced segments: when the interferer is a speech-shaped harmonic  
73 complex tone, SRTs below -10 dB have been reported for spatially collocated sound  
74 sources (Deroche and Culling, 2011). In contrast to these very low SRTs, observed in  
75 idealized laboratory conditions, Smeds et al. (2015) have presented evidence based on  
76 field recordings that, at least for hearing-aid users, real speech-to-noise ratios are  
77 rarely negative at all.

78         The present study is designed to create controlled virtual listening situations  
79 that are as realistic as possible, and to measure SRTs in those situations. At the same  
80 time, deviations from complete realism are included in order to access the relative  
81 roles of different perceptual mechanisms. To date, the most realistic simulations of  
82 this kind have been those of Culling (2013) and Westermann and Buchholz (2015),

83 and the present study shares features with each of these. However, unlike both these  
84 studies, the virtual room in Expt. 1 experimentally controls the presence of  
85 reverberation, while Expt. 2 is based on binaural room impulse responses (BRIRs)  
86 recorded from a real room, and so embodies all features of acoustic transmission,  
87 including the directivity of human speech production. In contrast to Culling (2013),  
88 but in common with Westermann and Buchholz, the masking sounds are continuous  
89 connected speech, as they would tend to be in a real listening situation. Compared to  
90 Westermann & Buchholz, the effect of the numerosity of the interferers is examined  
91 in greater detail (1, 2, 4 & 8, compared to 2 & 7), and reversed speech has been used  
92 as an additional form of masker. Among other things, these manipulations make it  
93 possible to discern the range of circumstances under which informational masking  
94 becomes apparent, and the SNRs at which normally hearing listeners can understand  
95 speech in realistic conditions.

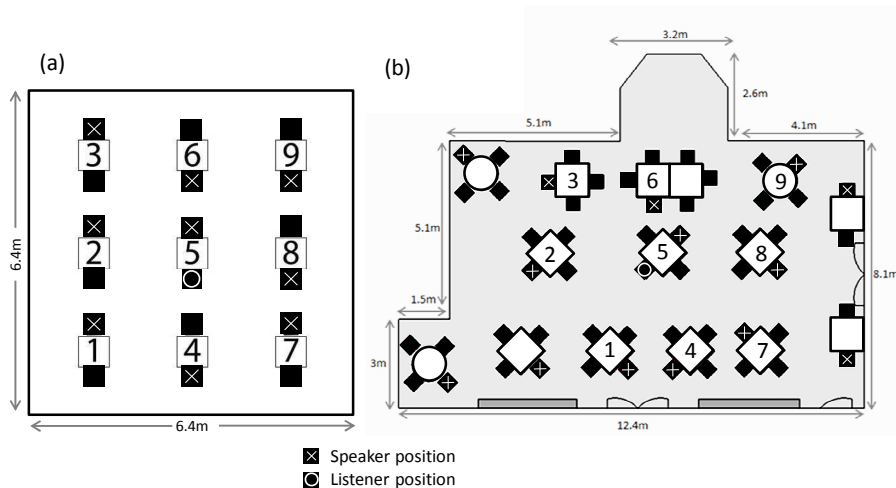
## 96 **II. METHODS**

97 The two experiments were similar in method except for the generation of the  
98 BRIRs and the spectral matching of target and interfering sources. In Expt. 1, BRIRs  
99 were generated by a ray-tracing algorithm as in Culling (2013), while in Expt. 2 they  
100 were recorded in a dining hall. In Expt. 1, the interfering speech was normalized, but  
101 was not matched to the target speech, while in Expt. 2, the target and interfering  
102 sources were filtered to match standardized speech spectra for the genders of the  
103 original recordings.

### 104 **A. BRIRs**

105 In Expt. 1, BRIRs for simulated restaurants, one *reverberant*, another *anechoic*,  
106 were generated using the image method of ray-tracing sound paths (Allen and  
107 Berkeley, 1979) and were identical to those of Culling (2013). For each sound path

108 between a source location and the listener's head, a head-related impulse response  
 109 (HRIR) was selected that was appropriate for that ray's angle of incidence with the  
 110 head. The HRIRs were recorded from a KEMAR by Gardner and Martin (1995). Each  
 111 was scaled and delayed according to the length and the surface interactions of the path  
 112 before being added into the combined BRIR. The restaurant was thus an empty box  
 113 with no furniture, sound sources were omnidirectional and surfaces reflected all  
 114 frequencies equally. Fig. 1a shows the layout, including the notional location of the  
 115 tables. The room was modelled to be 6.4 m square with a ceiling height of 2.5 m. In  
 116 the reverberant room, the surface absorbance of the floor, walls and ceiling were 0.07,  
 117 0.05 and 0.9, respectively. This gave a reverberation time ( $RT_{60}$ ) of 0.33 s. In the  
 118 anechoic room the absorbance was 1.0 for all surfaces. Source positions were  
 119 calculated on the basis that the room would contain nine regularly spaced tables for  
 120 two with the two people at each table 0.75 m apart. These BRIRs were 10,000  
 121 samples long at a sampling rate of 44.1 kHz (i.e. 227 ms in duration).



122

123

124

125

FIG.1. Table layouts used in each experiment. Left panel is a simulated restaurant with nine tables for two (Expt. 1). Right panel is Aberdare Hall at Cardiff University (Expt. 2).



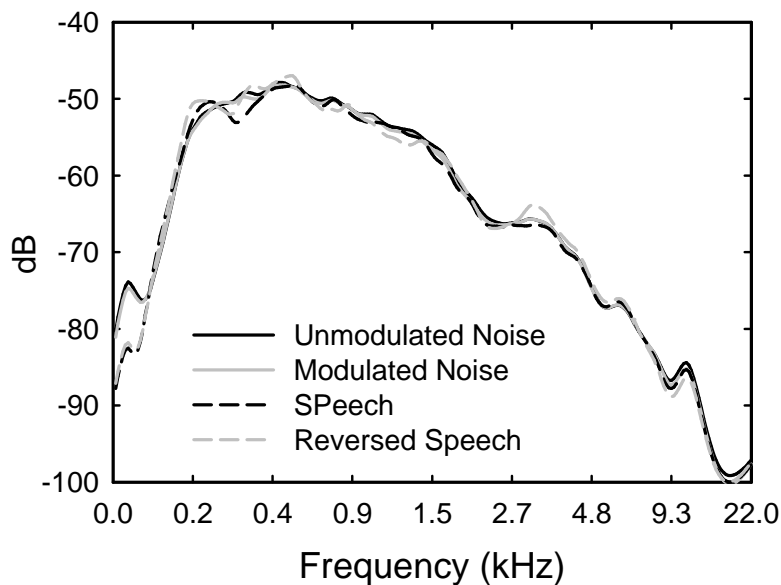
126 In Expt. 2, a *real* restaurant was used. BRIRs were recorded in Aberdare Hall at  
127 Cardiff University using the tone-sweep method (Müller and Massarini, 2001).  
128 Twenty-second logarithmic tone sweeps were presented from a B&K Head and Torso  
129 Simulator (type 4128), and recorded from a KEMAR manikin. The effect of  
130 KEMAR's ear canal resonance was removed from the BRIRs after recording by  
131 filtering them with a 512-point FIR filter designed to invert its diffuse field response,  
132 as measured by Killion (1979). Aberdare Hall can be divided in two by wooden  
133 panels. Recordings were made in the southern end of the hall with the dividing panels  
134 in place. This area is carpeted and partially wood-paneled, has approximate  
135 dimensions (L×W×H) of 12.4 m × 8.1 m × 4.5 m, and RT<sub>60</sub> of almost exactly 1  
136 second. It contains 14 tables for between 2 and 6 people (Fig. 1b). A speaker seat was  
137 selected at random for each table and BRIRs recorded between all selected speaker  
138 positions and a single listener position on the centrally located table 5. These BRIRs  
139 were 44,100 samples long (i.e. 1 second in duration).

## 140 **B. Interferers**

141 Recordings of monologues produced by four males and four females with a  
142 variety of British-English accents were selected from librivox audiobook recordings  
143 (librivox.org). Six-minute samples were drawn for each interferer. For the voices of  
144 each sex, the long-term excitation patterns (Moore and Glasberg, 1983) were  
145 equalized using specifically designed 512-point FIR filters. In Expt. 1 the interfering  
146 voices were equalized to each other using one of each sex as a model, but in Expt. 2  
147 they were equalized to published norms for male and female speech (Byrne et al.,  
148 1994, Table II). The rms power was also equalized. These speech interferers (SP)  
149 were then used to generate three other types of interferer, reversed speech (RS),  
150 speech modulated speech-shaped noise (MN) and unmodulated speech-shaped noise

151 (UN). Speech-shaping was achieved using a 512-point FIR filter designed to match  
 152 the long-term excitation pattern of either the male or female speech. Speech  
 153 modulation was achieved by extracting the modulation envelope through full-wave  
 154 rectification and low-pass filtering using a 512-point FIR filter with a 50 Hz cut-off.

155 The interferers were convolved with the BRIRs such that they were placed on  
 156 each of eight tables surrounding the listening position and then added together to  
 157 simulate different numbers of concurrent voices. The levels of the individual maskers  
 158 were attenuated by 3, 6 or 9 dB in order to compensate for the combination of 2, 4 or  
 159 8 interferers and keep the overall level of the masking complex constant. The  
 160 arrangement for each room is illustrated in Fig. 1, and the 5 distributions of voices in  
 161 the different conditions, which was designed to be similar across the two experiments,  
 162 is summarized in Table I.



163

164 FIG. 2. Long-term excitation patterns, based on 10 seconds of

165 material, of the four different types of interferer.

166

Interferers	Male	Female
1 male	3	
1 female		3
2	3	7
4	3, 9	1, 7
8	2, 3, 4, 9	1, 6, 7, 8

167 TABLE I. Table numbers selected for each number of interferers and  
 168 the genders of the voices (or noise spectra) placed on those tables.

169 Once the interferers were assembled, the excitation patterns (Moore and  
 170 Glasberg, 1983) were calculated in order to verify that each interferer type had the  
 171 same long-term masking potential. Example excitation patterns for the interferers  
 172 from Expt. 1 at the left ear and in the presence of 8 simultaneous interferers of each  
 173 type are plotted in Figure 2.

### 174 C. Targets

175 The target speech consisted of sentences from the IEEE corpus (Rothausser et  
 176 al. 1969), spoken by voice “DA” with an American-English accent. In Expt. 2 the  
 177 targets were, like the interferers, filtered to conform to Table II of Byrne et al. (1994).  
 178 These recordings were convolved with BRIRs for a speaker on the same table as the  
 179 listener (Table 5).

### 180 D. Procedure

181 Twelve participants with no known hearing impairments were recruited from  
 182 the Cardiff University undergraduate population for each experiment. They received  
 183 either payment or course credit for their participation. Participants were tested  
 184 individually in a single-walled audiometric booth with an auxiliary monitor visible

185 through the window for instructions and feedback. A keyboard inside the booth was  
186 provided for the participant to enter transcripts.

187 Expt.1 was run over two 90-minute sessions, while Expt. 2 was a single 90-  
188 minute session. Average completion time for each session was approximately 75  
189 minutes. Each experiment began with a detailed explanation of the SRT measurement  
190 procedure and a practice of the procedure. The practice consisted of two SRT  
191 measurements, one with two speech interferers and the other with two noise  
192 interferers. The spatial configurations employed differed from those used in the main  
193 experiment, consisting of two positions used only in the 8-interferer conditions.

194 In the experiments, the speech materials were presented in a fixed order while  
195 the experimental conditions were placed in a new, randomly generated sequence for  
196 each participant. For Expt. 1 there were 40 conditions, composed of 2 rooms  
197 (anechoic and reverberant), 5 interferer configurations (Table I), and 4 interferer types  
198 (SP, RS, MN & UN). In Expt. 2, there were only 20 conditions, because there was  
199 only one room.

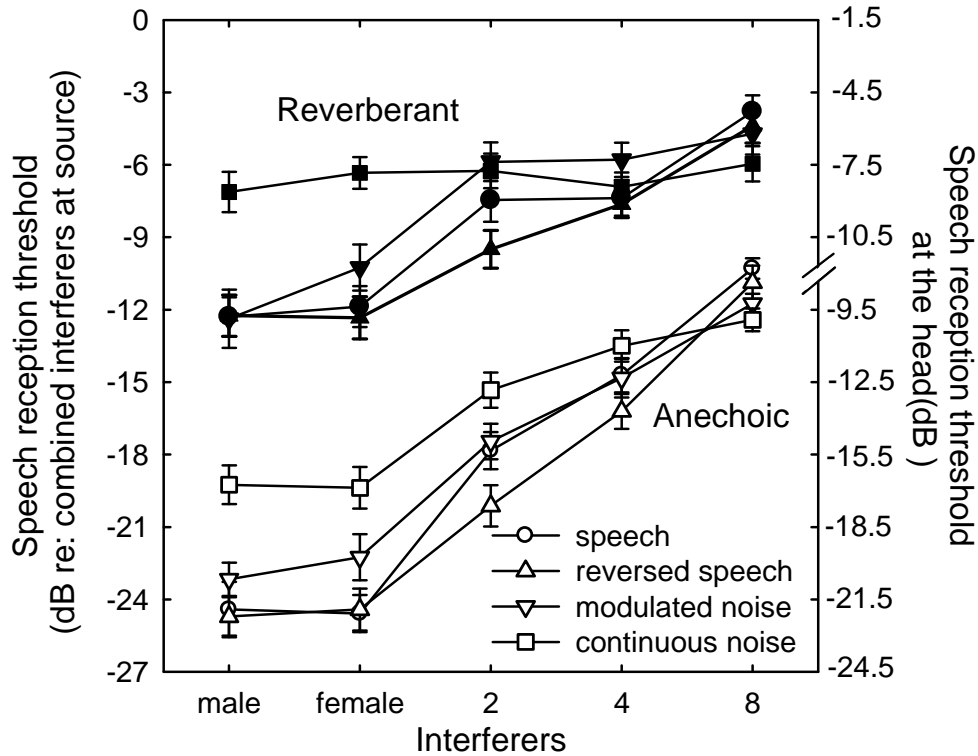
200 SRTs were measured using an adapted version of the Plomp and Mimpen  
201 (1979) method. The interfering sound started first and the participant initiated the first  
202 target sentence with a keypress. Participants listened for target sentences that were  
203 presented when “Listen for the target sentence” appeared on the auxiliary monitor.  
204 The speech-to-noise-ratio (SNR) was initially very low; the participant was instructed  
205 to press the enter key if they could not hear any of the first sentence. The sentence  
206 was repeated at a sound level that was 4 dB higher each time this was done. The  
207 participant was made aware that only two keywords correct would be needed to start  
208 the adaptive track. When the first transcript was entered, the words were checked  
209 automatically using a simple character-for-character match with the five keywords of

210 the stored transcript. If fewer than two words were correct, the participant was  
211 informed and the sound level of the first sentence was again increased by 4 dB. If at  
212 least two words were correct, the participant was then shown the actual transcript,  
213 with the five keywords in capitals and invited to self-score the transcript. The self-  
214 scoring method allows the participant to compensate for mis-typed and mis-spelled  
215 words as well as use of alternative spellings and homophones. Feedback on self-  
216 marking was provided by the experimenter after the practice. Once the two-word  
217 threshold was reached, the one-up/one-down adaptive track would begin. Each  
218 subsequent sentence was presented only once, participants did all their own marking  
219 and the sound level of the target speech was increased by 2 dB if the listener correctly  
220 identified less than 3 words. Otherwise the level was reduced by 2 dB. The entire  
221 interaction was recorded in detail in a log file in order to verify compliance with the  
222 instructions. Once all ten sentences in a list had been presented, the interfering sound  
223 was halted and the presentation levels that had been calculated after the last 8 trials  
224 was averaged to produce an estimate of the SRT.

### 225 **III. Results.**

226 Results from Expts 1 and 2 are shown in Figs. 3 and 4. The left ordinate  
227 indicates target speech levels at source compared to the total noise level at source.  
228 This measure does not reflect the SNR at the ear, because the target source is closer  
229 than the interferers. The right ordinates were therefore shifted to reflect the SNR of  
230 target speech against the interfering complex at the ear. The shift was calculated for  
231 the case of eight noise sources in order to minimize influence of interaural differences  
232 in interferer level. These SNRs were calculated using SII-weighted spectra (ANSI,  
233 1997) in order to compensate for spectral differences between the target and

234 interfering speech at source (in Expt. 1), and also differences in those spectra induced  
 235 by the room.

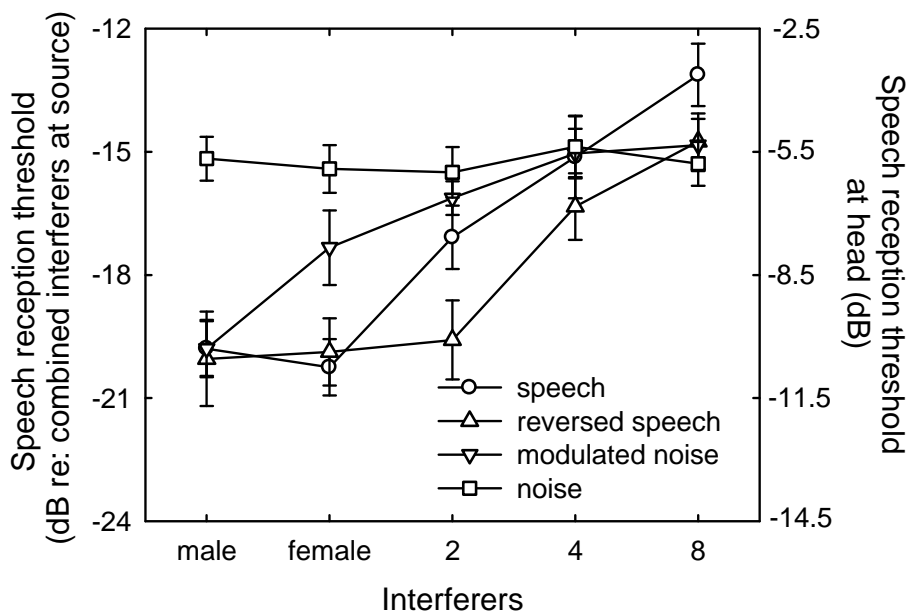


236

237 FIG. 3. Results from experiment 1. Speech reception thresholds for a voice  
 238 on the same table, as a function of the number/gender of interfering sources  
 239 at other tables. The ordinate indicates the signal-to-noise ratio at threshold  
 240 calculated on the basis of the source levels (i.e. before convolution with the  
 241 BRIRs). Filled symbols are for a simulated reverberant restaurant. Open  
 242 symbols are for a simulated anechoic restaurant. The right ordinate indicates  
 243 the approximate signal-to-noise ratio at the listener's head, based on the 8-  
 244 interferer condition. The right ordinate contains a break because the  
 245 introduction of reverberation reduces the signal-to-noise ratio at the head.  
 246 The upper section of the right ordinate thus applies to the reverberant  
 247 condition only and the lower section to the anechoic condition only.

248

249 The effects shown in Figs. 3 and 4 are reported here with respect to their  
 250 emergence in the statistical analysis. Each dataset was subjected to an analysis of  
 251 variance (ANOVA) with the factors room (anechoic vs. reverberant in Expt 1 only),  
 252 type of interferer (SP, RS, MN, UN) and number/gender of interferers (1 male, 1  
 253 female, 2, 4 and 8). Tukey HSD pairwise comparisons were used for post-hoc  
 254 analyses.



255

256 FIG. 4. Results from experiment 2. Speech reception thresholds for a  
 257 voice on the same table as a function of the number/gender of  
 258 interfering sources at other tables. The left ordinate indicates the  
 259 signal-to-noise ratio at threshold calculated on the basis of the source  
 260 levels (i.e. before convolution with the BRIRs). The right ordinate  
 261 indicates the approximate signal-to-noise ratio at the listener's head,  
 262 based on the 8-interferer condition.

263           The ANOVA for Expt. 1 revealed a significant main effect of room  
264 [ $F(1,11)=908, p<0.001$ ], reflecting higher SRTs in the reverberant room. There was  
265 also a significant effect of interferer type [ $F(3,33)=8.2, p<0.001$ ], reflecting a  
266 hierarchy among the interferers, in which continuous noise was the most effective  
267 interferer and speech and reversed speech were the least effective. All pairwise  
268 comparisons of interferer types were significant ( $p<0.01$ ). The number/gender of  
269 interferers also affected SRTs [ $F(4,44)=214, p<0.001$ ]; the SRTs increased  
270 significantly ( $p<0.01$ ) each time more voices were added, but SRT for one male or  
271 one female voice did not differ significantly. There was an interaction between the  
272 room and the number/gender of interferers [ $F(4,44)=44, p<0.001$ ], because the SRTs  
273 increased less steeply with the number of interferers in the reverberant room (see Fig,  
274 3). There was also an interaction between the type and number of interferers  
275 [ $F(12,132)=16.2, p<0.001$ ], in which the number of interferers had less effect for  
276 continuous noise than for the three modulated forms of interferer. No other  
277 interactions were significant.

278           The ANOVA for Expt. 2 revealed a very similar pattern for the real room with  
279 significant main effects of interferer type [ $F(3,33)=12.9, p<0.001$ ] and interferer  
280 number/gender [ $F(4,44)=37.0, p<0.001$ ], and a significant interaction between the two  
281 [ $F(12,132)=7.7, p<0.001$ ]. However, pairwise comparisons produced fewer  
282 significant differences. There were no longer significant differences between speech  
283 and reversed speech or between speech and modulated noise. Pairwise comparisons  
284 between different numbers of interferers no longer showed significant differences  
285 between a single female voice and a two-voice interferer ( $p=0.066$ ) and 4 and 8 voice  
286 interferers no longer differed significantly.



287 Pairwise comparison between different interferer types for the three different  
 288 rooms (the simulated anechoic and reverberant rooms from Expt. 1 and the real room  
 289 from Expt. 2) are summarized in Table II. These showed that, for the most part, the  
 290 unmodulated noise differed from the other three interferer types for one or two  
 291 interferers. However, in Expt. 2, reversed speech produced significantly lower SRTs  
 292 than both forward speech and speech modulated noise when two interferers were  
 293 present.  
 294

Interferer	Anechoic Room			Reverberant Room			Real Room		
	(Expt. 1)			(Expt. 1)			(Expt. 2)		
Number/type	RS	MN	UN	RS	MN	UN	RS	MN	UN
1 (male)	SP		**			**			**
	RS		**			**			**
	MN		*			**			**
1 (female)	SP		**			**		*	**
	RS		**			**			**
	MN		*			*			
2 interferers	SP						*		
	RS		**		*			**	**
	MN								
4 interferers	SP								
	RS		*						
	MN								

295 Table II. Results of Tukey HSD pairwise comparisons between the different interferer  
 296 types in the different rooms for each number of interferers (\* =  $p < 0.05$ ; \*\* =  $p < 0.01$ ).

## 297 **IV. Discussion**

298           The main objectives of the present study were to establish the role played by  
299 informational masking in realistic listening situations and to determine the lowest  
300 SNRs that can be tolerated by normally hearing listeners in such circumstances.  
301 Aspects of the data that are relevant to these two questions will therefore be addressed  
302 first.

### 303 **A. Informational masking.**

304           The role of informational masking in a realistic situation and normally hearing  
305 listeners was previously investigated by Westermann and Buchholz (2015). They  
306 concluded that the informational masking played a very limited role. This conclusion  
307 was based on the comparison of SRTs for speech interferers and unintelligible noise-  
308 vocoded interferers. The vocoded interferers were intended to produce the same  
309 amount of energetic masking as the speech interferers, including any benefits from  
310 modulation. The modulated speech-shaped noise interferers in the present experiment  
311 performed a similar role. Any addition of informational masking, produced by the  
312 speech, therefore could be observed as a relatively elevated SRT for speech  
313 interferers. A possible objection to this measure is that some release of masking will  
314 likely occur as a result of the harmonicity of the speech interferers (Deroche and  
315 Culling, 2011), an effect that would selectively lower the SRTs for speech interferers  
316 and so produce an underestimate of the informational masking effect.

317           In order to counter this objection, the present experiment also included  
318 reversed-speech interferers. Since these are unintelligible, but retain both modulation  
319 and harmonicity, they may provide a better baseline measure of energetic masking.  
320 Westermann and Buchholz did not observe elevated SRTs for speech interferers,  
321 compared to vocoded interferers, when the speech interferer was a different voice

322 from the target and was spatially separated from it (the more realistic case). The  
323 present data, however, do show some influence of informational masking with spatial  
324 separation. In most cases, the speech and reversed-speech interferers both provide the  
325 *lowest* SRTs, reflecting the benefits of modulation and harmonicity, but when there  
326 were two and perhaps four interferers, the reversed-speech interferer provided lower  
327 SRTs than the forward speech. This difference appears to reflect informational  
328 masking, presumably a specifically linguistic interference effect in which the listener  
329 is distracted by more than one intelligible interferer. The effect is more robust with  
330 two interferers with a difference apparent for all three rooms and reaching statistical  
331 significance in the case of the real room (Fig. 4). With four interferers, the mean SRT  
332 for reversed speech is lower than the others interferer types only in the case of an  
333 anechoic room, and this difference is non-significant. It seems likely that linguistic  
334 interference is already weak with four interferers and disappears in the presence of  
335 reverberation because reverberation impairs the intelligibility of the individual voices.  
336 These results are consistent with those previously found by Hawley et al. (2004).  
337 They observed higher SRTs from forward speech than reversed speech in anechoic  
338 conditions when there were two or three interferers, but not when there was only one.

339         The present study thus confirms, but qualifies Westermann and Buchholz's  
340 conclusions. It appears that a limited informational masking effect can be observed in  
341 realistic listening conditions, but only where there are a small number of interferers. It  
342 is also possible that further improvements to the stimuli might yet reveal a more  
343 extended role. There are two considerations, here.

344         First, although the use of reversed speech emulates the benefits of modulation  
345 and harmonicity in normal speech maskers, it may, at the same time, retain some  
346 informational masking potential. Hawley et al. (2004) noted that both reversed- and

347 forward-speech interferers seemed to facilitate an enhanced effect of spatial release  
348 from masking (by 2-3 dB) compared to interferers based on noise. The enhanced  
349 effect occurred for two or three interferers, but not when there was only one. They  
350 interpreted this result as a release from informational masking, which implies that  
351 *both* forward and reversed speech were generating informational masking when they  
352 were collocated with the target. Hawley et al. suggested that reversed speech may  
353 generate interference at lower levels of linguistic processing, such that, while it may  
354 not lead to intruding words or phrases, reversed speech might confuse mechanism of  
355 phonetic analysis. One approach to improving the emulation of energetic masking  
356 might be to use a speech-modulated complex tone, such that it possesses modulation  
357 and harmonicity, but no phonetic cues.

358         Second, the spatial set-up of the experiment placed all interferers roughly  
359 equidistant from the listener. Although this is a plausible configuration and makes a  
360 neat experimental design, many other real-life situations would have interferers at a  
361 variety of distances. In that case, those closer to the listener would tend to stand out  
362 and may have greater potential to induce informational masking.

### 363 **B. Real-life SNRs.**

364         The SNRs experienced and tolerated by people in the real world are essentially  
365 unknown, making it difficult to design appropriate signal processing for hearing aids  
366 or to generate acoustic standards for rooms. For instance, Rindel (2012) assumed that  
367 the lowest tolerable SNR in a room would be -3 dB on the basis that this is the  
368 approximate SRT for normally hearing listeners in continuous diffuse noise, but this  
369 assumption neglects, among other things, the possibility that the noise is more  
370 structured.

371 In order to address the absence of empirical data, Smeds et al. (2015) recorded  
372 the everyday acoustic exposure of 20 hearing-aid users for a total of 28 hours using  
373 bilateral microphones. Researchers analyzed these recordings, extracting segments  
374 containing speech addressed to the hearing-aid user and contemporaneous segments  
375 of background noise. A calculation was then made to obtain the SNR at which the  
376 speech had been received. The most striking result was that SNRs tended to be +5 dB  
377 or greater, suggesting that the frequent discussion of negative SNRs in the literature  
378 may be misguided. There are, however, a number of caveats that one should consider  
379 with respect to this finding.

380 First, the hearing aid users may have had strategies and habits that avoid  
381 exposure to poor SNRs, or friends and relations who seek to accommodate their  
382 difficulties by speaking loudly or during pauses in the noise. The reported SNRs may  
383 thus reflect the actual SNRs experienced by hearing-aid users during successful verbal  
384 interactions, but not the SNRs that they might like to be able to tolerate, nor the SNRs  
385 to which normally hearing listeners habitually expose themselves. Second, the method  
386 of deriving SNRs relies on the researcher correctly identifying acoustic segments  
387 when speech is addressed to the hearing-aid user, based only on listening to the  
388 recorded sound. It may be that segments at lower SNRs were more difficult to  
389 identify, and are consequently under-represented in the data. Finally, the hearing aid  
390 users were (unavoidably) placed in control of the recording process and may have  
391 biased their sampling of the acoustic environment in some way.

392 The present experiment, and that of Buchholz and Westermann (2015), took a  
393 completely different approach, in which we attempted to bring the real-world into the  
394 laboratory. In the present study, very realistic listening situations were created, and  
395 then the SRTs for 50% intelligibility of IEEE sentences were measured. The approach

396 has a number of limitations. It assumes that, in the real world, listeners will regularly  
397 place themselves in situations in which they can only just cope, so that measuring the  
398 threshold of coping informs us about real-life SNRs. The assumption is based upon  
399 the anecdotal experience that difficult listening situations, while not being prevalent,  
400 are sufficiently commonplace to be interesting. It also assumes that 50% intelligibility  
401 of standard sentence corpora occurs at a similar SNR to understanding well enough to  
402 sustain a real conversation. IEEE sentences are rather unpredictable compared to  
403 conversational speech, decreasing their intelligibility, but on the other hand, they are  
404 very clearly articulated. Greater than 50 % intelligibility is probably needed for  
405 conversation. Finally, the stimuli are also audio-only, and in real life one may expect  
406 SRTs to be improved by several dB by the use of lip-reading (Macleod and  
407 Summerfield, 1987). In order to address these limitations, a more realistic listening  
408 task will be required.

409         Notwithstanding these limitations, SRTs were found to increase with  
410 increasing numbers of interferers, even though the levels of individual interferers  
411 were adjusted in order to compensate for the increased masking energy. The increase  
412 in SRT was therefore attributable to the progressive degradation of perceptual  
413 unmasking mechanisms. We can thus see that the lowest tolerable SNR is  
414 considerably dependent upon the complexity of the listening scene. Because the effect  
415 of the number of interferers on overall sound level was compensated, the level of a  
416 given interferer reduces as the number of interferers increases. For a single interferer,  
417 an SRT of 0 dB (from the left ordinate) would thus represent a situation in which the  
418 interferer was speaking with the same effort as the target voice, but for 2, 4 and 8  
419 interferers, the SRT at this point would be -3, -6 and -9 dB, respectively. Bearing this  
420 in mind, we can see that only in the simulated reverberant restaurant with 2 or more

421 interferers (Expt. 1) does the target voice need to be raised above the level of the  
422 interfering voices in order to be heard; the real dining hall (Expt. 2) was thus a  
423 relatively benign environment with up to 8 interferers.

424         In a real listening environment, the background noise level will increase with  
425 increasing room occupancy, and the increase will be accentuated by the Lombard  
426 effect, an involuntary increase in vocal output induced by background noise (Lane and  
427 Tranel, 1971). This increase in vocal output is less than the increase in noise level,  
428 but, assuming that it is evenly distributed, will not change SNRs. However, once  
429 speech becomes unintelligible when produced at the same level as the interfering  
430 voices, as occurred in the reverberant room of Expt. 1, the various speakers in the  
431 room will come into direct competition. In Rindel's (2012) terms, the "acoustic  
432 capacity" of the room has been exceeded. This will make communication very  
433 difficult, and may induce a more marked increase in noise level (Maclean, 1959) or a  
434 behavioral adjustments such as leaning forward, or head orientation (Grange and  
435 Culling, 2016).

436         In order to compare with conventional SRT measurements without room  
437 simulations, the SRT at the head is indicated on the right ordinate in Figs. 3 and 4. We  
438 can see that in a simple scene with only one interferer, such as trying to hear what  
439 someone else is saying when the radio is on or against the noise of a vacuum cleaner,  
440 listeners can manage, in moderate reverberation (Fig. 4), at -5 to -10 dB SNR  
441 depending on the nature of the source, but as the scene becomes more complex SNRs  
442 need to be higher. Nonetheless, the most complex scenes examined here still produced  
443 SRTs approaching -5 dB, somewhat lower than the -3 dB assumed by Rindel (2012).

**444 C. Effects of reverberation.**

445 SRTs were lowest in the anechoic room, higher in the real room ( $RT_{60} = 1$  s)  
446 and highest in the simulated reverberant room ( $RT_{60} = 0.33$  s). The differences in SRT  
447 mainly reflect the detrimental effect of reverberation on mechanisms for perceptual  
448 separation. Reverberation reduces and distorts binaural differences generated by the  
449 interfering sound, and so affects spatial release from masking (Plomp, 1976;  
450 Lavandier and Culling, 2007, 2008). Reverberation distorts the harmonicity of  
451 interfering sounds when the fundamental frequency changes over time, leading to less  
452 effective harmonic cancellation (de Cheveigné, 1998; Culling et al., 2003; Deroche  
453 and Culling, 2011). Reverberation also temporally smears the masking sound such  
454 that temporal dips are filled in (Colin & Lavandier, 2013), and smears the target  
455 speech so that it becomes less intelligible (Houtgast and Steeneken, 1985). However,  
456 the detrimental effects of reverberation on unmasking from the interfering sound  
457 occur at lower levels if reverberation than the influences on temporal smearing of the  
458 target speech (Lavandier and Culling, 2008; Deroche and Culling, 2011).

459 It is noteworthy that the room with the highest  $RT_{60}$  was not the room with the  
460 highest SRTs. Beutelmann and Brand (2006) previously observed that spatial release  
461 from masking was not ordinally related to the  $RT_{60}$  of different rooms. Indeed,  
462 Culling et al. (2013) have argued that  $RT_{60}$  is a completely inappropriate statistic for  
463 considering speech intelligibility in noise, particularly if its interpretation is not  
464 moderated by room volume and likely source distances. In general, the direct-to-  
465 reverberant ratio of the interferers is a more accurate guide to the influence of  
466 reverberation. The direct-to-reverberant ratio is a statistic linked to the particular  
467 configuration of the source and receiver locations in the room, and so cannot be used  
468 to describe the room itself, but only a particular listening situation.



469           The increase in SRT with increasing numbers of interferers was also  
470 moderated by room reverberation. As more reverberation and more sources are added  
471 each situation approaches a completely diffuse continuous noise, as assumed by  
472 Rindel (2012). The slope of this increase in SRT with number of interferers is  
473 therefore strongly influenced by the starting SRT. If perceptual separation of the  
474 target and interfering noise is very good with a single interferer, then there is more  
475 separation effect to lose when the listening situation is made more complex.

#### 476 **D. Ever greater realism.**

477           In general, any area in which realism is limited leaves a study open to the  
478 criticism that results from the laboratory cannot be generalized. Both Westermann and  
479 Buchholz and the current experiments have moved to the use of continuous interfering  
480 sound, based on extended speech recordings. Preparation and presentation of such  
481 material is not as challenging as it once was. It is unclear whether this made much  
482 difference to the results obtained, but it certainly makes a difference to the realism  
483 experienced by the participants, who had a strong sensation of being immersed in the  
484 simulated environment. The technique saves the experimenter from any concerns  
485 about artefacts produced by the relative gating of the target and interferer, such as  
486 simultaneous sentence onsets being unusually confusing.

487           As noted above, the target speech was less realistic. In order to address the  
488 differences between listening to standardized speech corpora and real conversation,  
489 the most obvious route is to introduce real verbal interactions. Some work with real  
490 verbal interaction in noise has been pioneered by Cooke and Lu (2010), albeit in the  
491 context of studying speech production in these circumstances. Cooke and Lu had  
492 participants engage in conversation in order to solve a Sudoku puzzle together. In  
493 order for the technique to be adapted for use in an intelligibility measurement, the

494 speech level delivered from one interlocutor to the other will either need to be  
495 controlled, or monitored. While monitoring the level will place it under the control of  
496 the speaker, one may expect that the speaker will adapt it to a sufficient level to  
497 sustain the conversation, and this might make a reasonable outcome measure.

498         Westermann and Buchholz, used a commercial program, ODEON (Rindel,  
499 2000) to generate their BRIRs. This program enabled them to include furniture,  
500 frequency-dependent surface reflections and variations in reflectance across a given  
501 surface (e.g. windows within walls), but sound sources would still have been  
502 omnidirectional. The scene was then rendered over a loudspeaker array, which  
503 allowed listeners to make head movements, if desired, and to hear appropriate  
504 changes to the sound. Expt. 2 of the present study used real-room BRIRs that did  
505 capture source directionality using the mouth simulator of a B&K HATS. The scene  
506 was then rendered over headphones, which did not allow appropriate changes to the  
507 sound with head rotation. Since head rotation away from the target source has been  
508 shown to improve SRTs in noise (Grange and Culling, 2016), it would seem desirable  
509 to be able to recreate this aspect of real listening, but since it might also introduce an  
510 uncontrolled element in the results it would also be desirable that head orientation be  
511 continuously monitored. This could be achieved by adding a head tracker to the  
512 arrangements used by Westermann and Buchholz, or by using a head tracker to  
513 appropriately modify the stimulus in headphone presentation. The latter approach  
514 could be realized by preparing multiple versions of the target and the interferer,  
515 appropriate to different head orientations, and cross-fading between them as the head  
516 is turned.

517         No study to date, has attempted to include visual information in a realistic  
518 listening simulation. At a basic level, this would be a fairly simple addition, since it

519 would only require video presentation of the target speaker's face on a screen. This  
520 change would introduce the effect of lip-reading. Effects of lip-reading on speech  
521 intelligibility in noise are well-known (e.g. Macleod and Summerfield, 1987), and can  
522 be substantial in both normally hearing and hearing-impaired listeners. The benefits of  
523 rendering a more complete visual scene are less obvious and would require  
524 considerably greater effort. Nonetheless, effects on performance of competition from  
525 "distracter" faces have been observed (Yi et al. 2013), suggesting that truly realistic  
526 results can only be obtained with audio-visually rendered interferers. In any case, a  
527 more complex presentation system will be needed in order to simulate social  
528 interactions that include an exchange of conversation between multiple individuals,  
529 rather than the classic case of simply trying to recover a single voice from noise.

#### 530 **IV Conclusions**

531 Realistic simulations of listening situations that would typically be  
532 experienced in a restaurant indicate the speech reception threshold varies greatly with  
533 the complexity of the listening situation. Simple cases (one interfering voice) permit  
534 SRTs of around as low as -10 dB, but more complex cases can elevate SRTs to -5 dB.  
535 Informational masking is observed in realistic listening conditions under quite limited  
536 conditions; in the present case, it was only observed when two interferers were  
537 present.

#### 538 **V Acknowledgments**

539 Sasha Priddy assisted in the collection of BRIRs during vacation scholarship  
540 funded by the School of Psychology. Jacques Grange assisted with experimental data  
541 collection. I am grateful for comments by two anonymous reviewers and for those of  
542 Mickael Deroche, Mathieu Lavandier and Adam Westermann on a presubmission

- 543 draft of the manuscript. Tony Watkins kindly loaned the KEMAR and the B&K  
544 HATS for BRIR collection.

## References

- Allen, J. B., Berkley, D. A., and Hill, M. (1979). "Image method for efficiently simulating small-room acoustics," *J. Acoust. Soc. Am.*, 65, 943–950.
- ANSI (1997). S3.5. Methods for the Calculation of the Speech Intelligibility Index. American National Standards Institute, New York.
- Beutelmann, R., and Brand, T. (2006). "Prediction of speech intelligibility in spatial noise and reverberation for normal-hearing and hearing-impaired listeners," *J. Acoust. Soc. Am.*, 120, 331–342.
- Bronkhorst, A. W., and Plomp, R. (1988). "The effect of head-induced interaural time and level differences on speech intelligibility in noise," *J. Acoust. Soc. Am.*, 83, 1508–1516.
- Brungart, D. S., Simpson, B. D., Ericson, M. a, and Scott, K. R. (2001). "Informational and energetic masking effects in the perception of multiple simultaneous talkers," *J. Acoust. Soc. Am.*, 110, 2527–2538.
- Byrne, D., Dillon, H., Tran, K., Arlinger, S., Wilbraham, K., Cox, R., Hagerman, B., Hetu, R., Kei, J., Lui, C., Kiessling, J., Nasser Kotby, M., Nasser, N. H. A., El Kholy, W. A. H., Nakanishi, Y., Oyer, H., Powell, R., Stephens, D., Meredith, R., Sirimanna, T., Tavartkiladze, G., Frolenkov, G. I., Westerman, S., and Ludvigsen, C. (1994). An international spectra comparison of long-term average speech. *J. Acoust. Soc. Am.*, 96, 2108–2120.
- Collin, B., & Lavandier, M. (2013). "Binaural speech intelligibility in rooms with variations in spatial location of sources and modulation depth of noise interferers." *J. Acoust. Soc. Am.*, 134, 1146–1159.

- Cooke, M., and Lu, Y. (2010). "Spectral and temporal changes to speech produced in the presence of energetic and informational maskers," *J. Acoust. Soc. Am.*, 128, 2059–2069.
- Culling, J. F. (2013). "Energetic and informational masking in a simulated restaurant environment" in Moore, B. C. J., Carlyon, R. P., and Gockel, H., Patterson, R. D. and Winter, I. M.. (eds) *Basic Aspects of Hearing: Physiology and Perception* (Springer, New York)
- Culling, J. F., Hodder, K. I., and Toh, C. Y. (2003). "Effects of reverberation on perceptual segregation of competing voices," *J. Acoust. Soc. Am.*, 114, 2871-2876.
- Culling, J. F. Lavandier, M. and Jelfs, S. (2013). "Predicting binaural speech intelligibility in architectural acoustics" in Blauert, J. (Ed.) *The Technology of Binaural Listening* (Springer, Heidelberg)
- de Cheveigné, A. (1998). "Cancellation model of pitch perception," *J. Acoust. Soc. Am.*, 103, 1261–1271.
- Deroche, M. L. D., and Culling, J. F. (2011). "Voice segregation by difference in fundamental frequency: evidence for harmonic cancellation," *J. Acoust. Soc. Am.*, 130, 2855–2865.
- Fletcher, H. and Galt, R. H. (1950). "The perception of speech and its relation to telephony," *J. Acoust. Soc. Am.*, 22, 89–151.
- Gardner, W. G., and Martin, K. D. (1995). "HRTF measurements of a KEMAR," *J. Acoust. Soc. Am.*, 97, 3907–3908.
- Grange, J. A., & Culling, J. F. (2016). The benefit of head orientation to speech intelligibility in noise. *J. Acoust. Soc. Am.*, 139, 703–712.

- Hawley, M. L., Litovsky, R. Y., and Culling, J. F. (2004). "The benefit of binaural hearing in a cocktail party: Effect of location and type of interferer," *J. Acoust. Soc. Am.*, 115, 833–843.
- Houtgast, T., and Steeneken, H. J. M. (1985). "A review of the MTF concept in room acoustics and its use for estimating speech intelligibility in auditoria," *J. Acoust. Soc. Am.*, 77, 1069–1077.
- Killion, M. C. (1979). "Equalization filter for eardrum-pressure recording using a KEMAR manikin," *J. Audio. Eng. Soc.*, 27, 13–16.
- Lane H. and Tranel B. (1971). "The Lombard sign and the role of hearing in speech" *J. Speech Hear Res* 14 (4): 677–709.
- Lavandier, M., and Culling, J. F. (2007). "Speech segregation in rooms: effects of reverberation on both target and interferer," *J. Acoust. Soc. Am.*, 122, 1713–1723.
- Lavandier, M., and Culling, J. F. (2008). "Speech segregation in rooms: monaural, binaural, and interacting effects of reverberation on target and interferer," *J. Acoust. Soc. Am.*, 123, 2237–2248.
- Licklider, J. C. R. (1948). "The influence of interaural phase relations upon the masking of speech by white noise," *J. Acoust. Soc. Am.*, 20, 150–159.
- Maclean, W. (1959). On the Acoustics of Cocktail Parties. *J Acoust Soc Am*, 31, 79–80.
- Macleod, A. and Summerfield, Q. (1987) "Quantifying the contribution of vision to speech perception in noise" *Br. J. Audiol.* 21, 131-142.
- Miller, G. A. (1947). "The masking of speech," *Psychol. Bull.*, 44, 105–129.

- Moore, B. C. J., and Glasberg, B. R. (1983). "Suggested formulae for calculating auditory-filter bandwidths and excitation patterns," *J. Acoust. Soc. Am.*, 74, 750–753.
- Müller, S. and Massarini, P. (2001). Transfer function measurement with sweeps, *J Audio Eng Soc*, 49, 443–471.
- Plomp, R. (1976). "Binaural and monaural speech intelligibility of connected discourse in reverberation as a function of azimuth of a single competing sound source (speech or noise)," *Acustica*, 34, 200–211.
- Plomp, R., and Mimpen, A. (1979). "Improving the reliability of testing the speech reception threshold for sentences," *Audiology*, 18, 43–52.
- Rhebergen, K. S., & Versfeld, N. J. (2005). A Speech Intelligibility Index-based approach to predict the speech reception threshold for sentences in fluctuating noise for normal-hearing listeners. *J Acoust. Soc. Am.*, 117, 21812–2192.
- Rindel, J. (2000). "The use of computer modeling in room acoustics," *J. Vibroengin.* 3, 219–224.
- Rindel, J. H. (2012). "Acoustical capacity as a means of noise control in eating establishments" Joint Baltic-Nordic Acoustics Meeting (Odense, Denmark).
- Rothausler, E. H., Chapman, W. D., Guttman, N., Hecker, M. H. L., Nordby, K. S., Silbiger, H. R., Urbanek, G. E., and Weinstock, M. (1969). "IEEE recommended practice for speech quality measurements." *IEEE Trans Aud. Electroacoust.*, 17, 227–246.
- Smeds, K., Wolters, F., and Rung, M. (2015). "Estimation of signal-to-noise ratios in realistic sound scenarios," *J. Am. Acad. Audiol.*, 26, 183–196.



Westermann, A., and Buchholz, J. M. (2015). "The effect of spatial separation in distance on the intelligibility of speech in rooms," *J. Acoust. Soc. Am.*, 137, 757–67.

Yi, A., Wong, W., and Eizenman, M. (2013). "Gaze patterns and audiovisual speech enhancement," *J. Speech. Lang. Hear. Res.*, 56, 471–80.