



HAL
open science

Investigating the role of harmonic cancellation in speech-on-speech masking

Luna Prud'homme, Mathieu Lavandier, Virginia Best

► **To cite this version:**

Luna Prud'homme, Mathieu Lavandier, Virginia Best. Investigating the role of harmonic cancellation in speech-on-speech masking. *Hearing Research*, 2022, 426, pp.108562. 10.1016/j.heares.2022.108562 . hal-03886240

HAL Id: hal-03886240

<https://hal.science/hal-03886240v1>

Submitted on 11 Dec 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Investigating the role of harmonic cancellation in
2 speech-on-speech masking

3 Luna Prud'homme^a, Mathieu Lavandier^{a,*}, Virginia Best^b

4 ^a*Univ. Lyon, ENTPE, Laboratoire de Tribologie et Dynamique des Systèmes UMR 5513,*
5 *Rue M. Audin, 69518 Vaulx-en-Velin Cedex, France*

6 ^b*Department of Speech, Language and Hearing Sciences, Boston University, 635*
7 *Commonwealth Ave, Boston, MA, 02215, USA*

8 **Abstract**

This study investigated the role of harmonic cancellation in the intelligibility of speech in “cocktail party” situations. While there is evidence that harmonic cancellation plays a role in the segregation of simple harmonic sounds based on fundamental frequency (F0), its utility for mixtures of speech containing non-stationary F0s and unvoiced segments is unclear. Here we focused on the energetic masking of speech targets caused by competing speech maskers. Speech reception thresholds were measured using seven maskers: speech-shaped noise, monotonized and intonated harmonic complexes, monotonized speech, noise-vocoded speech, reversed speech and natural speech. These maskers enabled an estimate of how the masking potential of speech is influenced by harmonic structure, amplitude modulation and variations in F0 over time. Measured speech reception thresholds were compared to the predictions of two computational models, with and without a harmonic cancellation component. Overall, the results suggest a minor role of harmonic cancellation in reducing energetic masking in speech mixtures.

9 *Keywords:* Speech Intelligibility, Harmonic Cancellation, Auditory Modeling,
10 Binaural Hearing

*Corresponding author

Email address: mathieu.lavandier@entpe.fr (Mathieu Lavandier)

11 1. Introduction

12 In “cocktail party” scenarios (Cherry, 1953), where speech must be under-
13 stood in the presence of noise and competing talkers, there are several factors
14 that may improve intelligibility, including spatial separation of sources, masker
15 amplitude modulation, and differences in harmonic structure between sources.
16 While the effects of spatial separation and masker amplitude modulation have
17 been extensively characterized and incorporated into computational models of
18 speech intelligibility (Lavandier & Best, 2020), the role of harmonicity is less
19 well understood.

20 In the context of speech-on-speech masking, several previous studies showed
21 a beneficial effect of fundamental frequency (F0) differences. Brokx & Noot-
22 boom (1982) tested the intelligibility of monotonized target speech against
23 monotonized masker speech. They found that the percentage of errors decreased
24 with increasing F0 difference between target and masker, except when the dif-
25 ference was one octave. Darwin et al. (2003) found that differences in F0 and
26 vocal tract length both improved segregation of target and masker talkers. In a
27 recent study, Popham et al. (2018) tested speech intelligibility with target and
28 masker that were either both natural (harmonic) speech or both inharmonic
29 speech. The percentage of correctly reported words decreased when the speech
30 was inharmonic, suggesting that harmonicity is important for the successful
31 segregation of speech mixtures. In each of these studies, it is not entirely clear
32 what mechanisms underlie the observed benefits. One of the complicating fac-
33 tors is that speech-on-speech situations involve both energetic masking (EM,
34 Culling, 2016) and informational masking (IM, Kidd et al., 2016). Whereas EM
35 refers to masking that renders target speech inaudible, IM refers to masking
36 that happens even when the target speech is audible, and is often attributed to
37 the listener’s inability to perceptually segregate competing voices and selectively
38 attend to the target talker. Harmonicity-based effects could theoretically reduce
39 both kinds of masking. For example, harmonic structures may be readily sup-
40 pressed by the auditory system, thus reducing EM. Alternatively, differences in

41 F0 may enhance the perceptual segregation of competing talkers and effectively
42 reduce IM.

43 The present study focuses on harmonicity-based effects on EM, and delib-
44 erately puts aside harmonicity-based effects on IM. A handful of recent studies
45 are highly relevant here as they measured speech intelligibility in the presence
46 of harmonic complexes (often called buzzes; Leclère et al., 2017; Steinmetzger
47 & Rosen, 2015; Deroche & Culling, 2013). These stimuli have the benefit that
48 they are harmonic, but they are not speech-like and thus they can be assumed
49 to cause primarily EM. These studies demonstrated that harmonic sounds exert
50 less EM than inharmonic sounds, and considered several possible mechanisms.
51 Deroche & Culling (2013) suggested that listeners might be able to extract tar-
52 get information in the spectral dips of harmonic maskers, between the resolved
53 partials (“spectral glimpsing”). Another possibility is that listeners may detect
54 the harmonic structure of the masker and suppress it when its F0 is different
55 from the one of the target (harmonic cancellation, de Cheveigné, 2021, 1993).
56 It is unclear though whether this reduction in EM translates to natural sounds
57 like speech. With speech maskers, there are several co-existing acoustic charac-
58 teristics to consider that may complicate the picture: the presence of unvoiced
59 segments, intonation, and amplitude modulation. Indeed, there are some in-
60 dications that these factors may reduce the benefits associated with harmonic
61 maskers (Deroche & Culling, 2013; Leclère et al., 2017; Deroche & Gracco, 2019).
62 For example, Leclère et al. (2017) showed that the advantage due to harmonicity
63 was reduced for intonated buzzes compared to monotonized buzzes. Deroche
64 & Gracco (2019) tested speech intelligibility against harmonic complexes and
65 monotonized speech maskers with one or two F0s. They found different results
66 for the two types of masker, suggesting that F0-based unmasking might operate
67 differently for speech and buzz maskers.

68 The current study aimed to extend these findings by comparing speech in-
69 telligibility for a range of harmonic and inharmonic maskers (including noise,
70 buzzes and speech) in the same group of listeners under otherwise identical
71 conditions. By comparing performance across masker types, the aim was to

72 estimate the influence of harmonicity while controlling for various other charac-
73 teristics of speech. IM was deliberately minimized in order to focus on energetic
74 aspects of harmonicity-based unmasking.

75 The present study is also part of a larger effort to develop a speech intelli-
76 gibility model that accurately captures the EM caused by complex interfering
77 sounds such as speech. Our motivation for doing so is to have a means by
78 which the relative contributions of EM and IM can be confidently estimated
79 in real-world listening situations. Previous studies have estimated EM caused
80 by speech maskers using models validated only for amplitude-modulated noise
81 maskers, thus implicitly (Biberger & Ewert, 2019) or explicitly (Lavandier et al.,
82 2021; Wasiuk et al., 2020) assuming that masker harmonicity was not relevant.
83 Steinmetzger et al. (2019) tested four models on speech intelligibility data col-
84 lected in the presence of speech-shaped noise and intonated harmonic complexes
85 (Steinmetzger & Rosen, 2015). All of the models underestimated the benefit due
86 to masker harmonicity. The authors speculated that the performance of the
87 models might be improved by implementing a mechanism of enhanced stream
88 segregation dependent on masker harmonicity. Prud’homme et al. (2020) pro-
89 posed a model with an implementation of harmonic cancellation that allowed
90 good predictions of intelligibility against stationary buzzes. This model was fur-
91 ther developed to take into account the detrimental effect of masker intonation
92 by Prud’homme et al. (2022), but was never tested in speech-on-speech sce-
93 narios. In the present study, by combining behavioral measurements and model
94 predictions, we aimed to investigate whether harmonicity (and harmonic cancel-
95 lation in particular) plays a role in reducing EM in speech-on-speech situations.
96 Specifically, we compared the predictions of models with and without harmonic
97 cancellation, to determine whether this mechanism is a necessary component of
98 a comprehensive speech intelligibility model.

99 **2. Main experiment**

100 *2.1. Methods*

101 *2.1.1. Listeners*

102 Nine listeners (ages 19-23 years, mean age 21) participated in the main ex-
103 periment (ten listeners originally participated but one listener had unusually
104 poor speech intelligibility, even in quiet, and thus their data were excluded from
105 the analysis). All listeners had normal pure tone thresholds (not exceeding 20
106 dB HL at octave frequencies from 250 to 8000 Hz) and were paid for their par-
107 ticipation. All procedures were approved by the Boston University Institutional
108 Review Board.

109 *2.1.2. Stimuli*

110 The target sentences were matrix sentences taken from a closed set (Kidd
111 et al., 2008). The sentences consisted of five words (name, verb, number, ad-
112 jective, object) and each word was drawn randomly from a set of eight options.
113 The target sentence was always spoken by the same North American female
114 voice (mean F0 = 180 Hz).

115 Seven types of masker were used: speech-shaped noise (SSN), monotonized
116 buzz, intonated buzz, monotonized speech, natural speech, vocoded speech
117 and reversed speech. SSN was used as a baseline without any harmonicity
118 or amplitude modulation. The monotonized and intonated buzzes are non-
119 speech maskers without the slow amplitude modulations that are characteristic
120 of speech, but with harmonicity (F0 fixed or varying over time, respectively).
121 Vocoded speech has amplitude modulations that are similar to speech, which
122 should provide a similar amount of temporal dip listening, but no harmonic-
123 ity. Monotonized speech was included because previous data suggested that
124 harmonic cancellation operates more effectively for monotonized than intonated
125 buzz maskers (Leclère et al., 2017). Reversed speech was added as a control
126 to confirm that our forward (natural) speech masker caused minimal IM: since
127 reversed speech causes very little IM then we expected no difference between
128 these two maskers.

129 The seven maskers were all derived from the same speech monologue spoken
130 by an Australian accented male talker (mean F0 = 112 Hz). The SSN was a
131 white noise filtered to have the same long-term excitation pattern (Glasberg
132 & Moore, 1990) as the monologue. The buzzes were harmonic complexes with
133 partials in random phase. The buzz was either monotonized with a fixed F0 of
134 112 Hz, or intonated with a continuous F0 contour extracted from the speech
135 monologue. This F0 contour was applied to the buzz using PRAAT PSOLA
136 (Boersma & Weenink, 2018). The monotonized speech was created by fixing
137 the F0 of the monologue to its mean (112 Hz) using PRAAT PSOLA. The
138 vocoded speech was created with an 8-channel vocoder, using an envelope low-
139 pass filter-cutoff frequency of 150 Hz. All maskers were passed through a 0.1-s
140 finite impulse response filter to match the average long-term excitation pattern
141 of the monologue. Figure 1 presents the excitation patterns of all maskers.

142 Stimuli were spatialized using anechoic KEMAR head related transfer-functions
143 (HRTFs; Gardner & Martin, 1995). The target was presented at 0° azimuth
144 and the masker was either presented at 0° azimuth (co-located condition) or
145 60° to the side (separated condition). This difference in location, as well as the
146 difference in talker sex, and the difference in the structure and content of the
147 speech materials, all served to distinguish the target from the masker and thus
148 minimize IM.

149 *2.1.3. Procedure*

150 The stimuli were presented via a 24-bit soundcard (RME HDSP 9632, Haimhausen,
151 Germany) to a pair of circumaural headphones (Sennheiser HD 280 Pro, Wede-
152 mark, Germany). The listeners were seated in a double-walled sound treated
153 booth. After listening to each masked sentence, they were presented with a grid
154 of 40 words (five categories x eight options). They were instructed to select one
155 word in each category and were then presented with correct answer feedback.

156 The experiment took two sessions of approximately two hours each. Each
157 session was composed of twenty-nine blocks of 20 trials. Within a block, the
158 masking condition was fixed but the signal-to-noise ratio (SNR) varied randomly

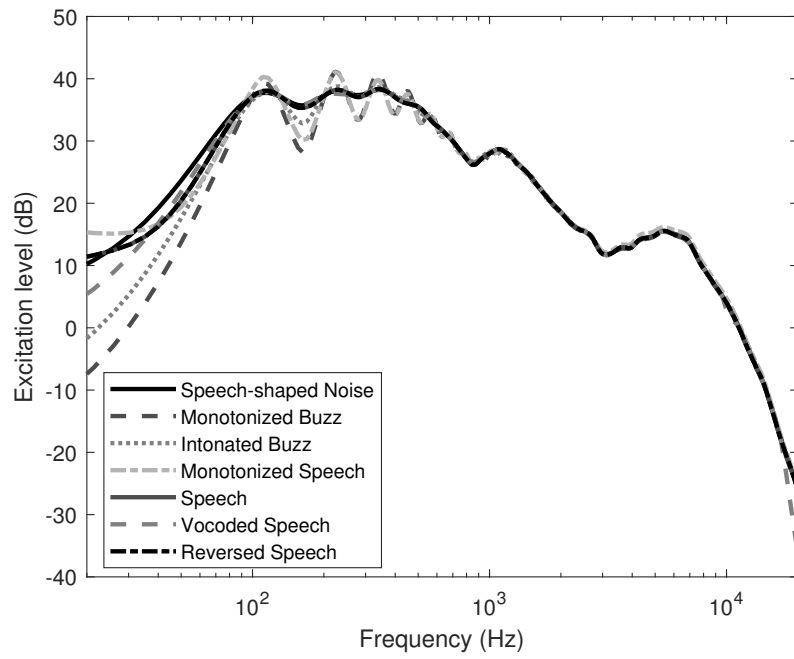


Figure 1: Long-term excitation patterns of the seven maskers tested in the main experiment

159 between five chosen values (from -40 to -10 dB). The masker presentation level
160 was fixed at 65 dB SPL and the target level was varied to achieve the required
161 SNR. The first block in each session consisted of target sentences alone (in quiet)
162 to familiarize the listeners with the task and to make sure that they were able to
163 understand the speech in quiet at the different target levels (which ranged from
164 25 to 55 dB SPL). This was followed by twenty-eight blocks of testing. Each
165 masking condition was presented twice per session. The order was randomized
166 across participants. At the end of the two sessions, each listener had performed
167 four blocks in each condition, which resulted in 80 scored words at each SNR.

168 The percentage of correct words was calculated for each participant at each
169 SNR in each condition. Logistic functions were fitted using the psignfit toolbox
170 version 4 for MATLAB, which implements the maximum likelihood method
171 described by Wichmann & Hill (2001). The lower asymptote was set to chance
172 performance (12.5%). Speech reception thresholds (SRTs) corresponding to 50%
173 correct were extracted from the logistic fits.

174 *2.2. Results*

175 All subjects (except the one whose data was excluded from the analysis)
176 performed well in quiet with scores above 85% correct at all target levels.

177 Figure 2, top panel, presents the mean SRTs across listeners for the different
178 masker types in the co-located and separated conditions. A repeated-measures
179 ANOVA was performed with two factors (masker type x spatial separation).
180 The main effects of spatial separation, masker type and their interaction were
181 significant ($F(1,8) = 940.99$, $p < 0.001$, $F(6,48) = 107.91$, $p < 0.001$ and $F(6,48) =$
182 15.56 , $p < 0.001$, respectively). Post hoc paired t-tests ($p < 0.05$) suggested that
183 in the co-located condition, SRTs were significantly higher for SSN than for the
184 buzzes, and were significantly higher for intonated than monotonized buzzes. A
185 similar pattern was found in the separated condition, with the exception that
186 the difference in SRT between SSN and intonated buzz was not significant. In
187 both spatial conditions, there was no significant difference in SRT between nat-
188 ural speech and the other speech-based maskers (monotonized speech, vocoded

189 speech, reversed speech). The important results here are that (a) the effect
190 of intonation observed for the buzz maskers was not replicated for the speech
191 maskers, (b) the absence of periodicity in the vocoded speech masker had little
192 effect on SRTs, and (c) forward and reversed speech maskers gave similar SRTs
193 consistent with minimal IM.

194 **3. Control experiments**

195 *3.1. Rationale*

196 SRTs in the main experiment were very low (-32 dB on average in the sepa-
197 rated conditions), presumably because of the large number of cues available to
198 reduce both EM and IM. Because of these low values, there was some concern
199 that a floor effect may have limited our ability to see differences across masker
200 types. Thus, two control experiments were conducted to increase the SRTs in
201 different ways. In control experiment 1, the speech task was made more difficult
202 by using open-set target sentences as opposed to closed-set matrix sentences. In
203 control experiment 2, in addition to using open-set materials, a second speech
204 masker was added to increase the overall amount of masking.

205 In addition to increasing the overall difficulty, the idea behind control exper-
206 iment 2 was that increasing the EM may increase the opportunities for harmonic
207 cancellation to operate. Specifically, by adding a second speech masker, there
208 should be more voiced parts and fewer dips overall in the speech masker. This
209 experiment was designed so that IM, which is known to be more prominent for
210 two-talker maskers than for single-talker maskers (Freyman et al., 2004; Brun-
211 gart et al., 2001), was still minimized as much as possible.

212 *3.2. Methods*

213 Five listeners (ages 19-22 years, mean age 21) who participated in the main
214 experiment also participated in control experiment 1. Five new listeners who did
215 not participate in the main experiment (ages 19-22 years, mean age 21; normal
216 pure tone thresholds at octave frequencies from 250 to 8000 Hz) participated in
217 control experiment 2.

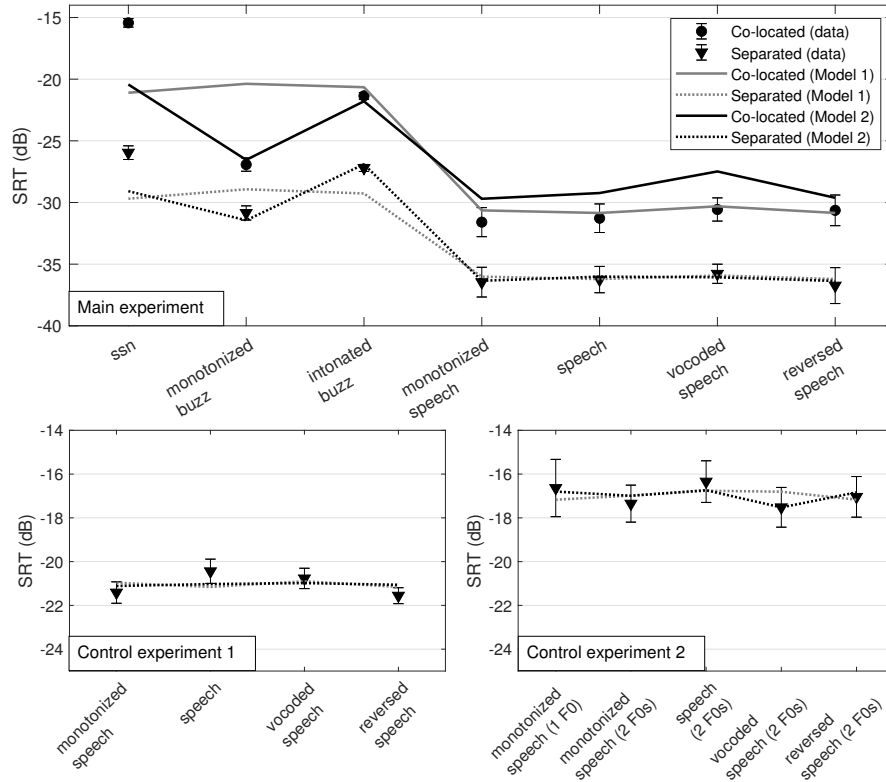


Figure 2: Mean SRTs with standard errors across participants measured in the main experiment (top panel), in the control experiment 1 (bottom left panel), and the control experiment 2 (bottom right panel) with the corresponding model predictions using: model 1, a binaural model without harmonic cancellation (Vicente & Lavandier, 2020) and model 2, a binaural model with harmonic cancellation (Prud'homme et al., 2022).

218 In the two control experiments, the target sentences were taken from the
219 Harvard Sentence List (Rothausser et al., 1969) and were composed of five key-
220 words. They were spoken by a female talker with a North American accent
221 (mean F0 = 190 Hz). Only the spatially separated configuration was tested,
222 and in control experiment 2, the two maskers were both presented at the same
223 location (60° azimuth).

224 In control experiment 1, only four masker conditions (those that could have
225 been influenced by a floor effect) were re-tested: speech, monotonized speech,
226 vocoded speech and reversed speech.

227 In control experiment 2, those same masker conditions were tested with two
228 maskers instead of one. The maskers were generated as in the main experiment,
229 derived from speech monologues spoken by male voices (one at a mean F0 of
230 112 Hz, same monologue as the main experiment, the other at a mean F0 of
231 130 Hz), both with an Australian accent. This resulted in four masker types:
232 two-talker natural speech, two-talker monotonized speech, two-talker vocoded
233 speech, two-talker reversed speech. A fifth masker type was constructed by
234 adding two time-shifted copies of the monotonized male monologue from the
235 main experiment. This resulted in a two-talker masker with a single steady F0
236 at 112 Hz. This masker was theoretically optimized for harmonic cancellation,
237 having a higher proportion of voiced energy than the single-talker version and
238 a single, steady F0.

239 In the control experiments, the participants were instructed to type the sen-
240 tence they heard. The correct transcript was then displayed on the screen, with
241 the keywords in capital letters and the participants had to self-mark their num-
242 ber of correct keywords. For each control experiment, the listeners performed
243 two sessions of one hour. Each session was composed of 20 blocks of ten sen-
244 tences in one of the conditions and at one SNR. Five SNRs were tested from
245 -30 to -10 dB. The listeners also performed a block in quiet at the start of each
246 session. SRTs were extracted as per the main experiment.

247 *3.3. Results*

248 Figure 2 presents the SRTs measured in control experiments 1 and 2, in
249 the bottom left and right panels, respectively. The SRTs in the control ex-
250 periments were higher than in the main experiment as intended. Repeated
251 measures ANOVAs found no significant effect of masker type in control experi-
252 ment 1 [$F(3,12) = 1.54$, $p = 0.26$] or control experiment 2 [$F(4,16) = 1.79$, $p =$
253 0.18]. These results corroborate the key results of the main experiment, again
254 suggesting that the harmonic structure of a speech masker has little impact on
255 SRTs under conditions of minimal IM.

256 **4. Modeling**

257 *4.1. Rationale*

258 In order to further investigate the potential role of harmonic cancellation in
259 the present study, two speech intelligibility models were applied to the stimuli.
260 The chosen models were the binaural model proposed by Vicente & Lavandier
261 (2020) and validated only for amplitude-modulated noise maskers, and the bin-
262 aural harmonic cancellation model proposed by Prud’homme et al. (2022) vali-
263 dated only for harmonic complex maskers (with both intonation and amplitude
264 modulation). By comparing how well each model can account for the data, we
265 aimed to provide further support for (or against) a role of harmonic cancellation
266 in reducing EM for speech maskers.

267 *4.2. Models*

268 The models tested here have the same structure. The inputs for both models
269 are the target and masker signals at the ears of the listener. The target input
270 is represented by an averaged signal obtained by adding several target sen-
271 tences (Vicente & Lavandier, 2020). The masker signal is segmented into time
272 frames using half-overlapping Hann windows. The signals are passed through a
273 Gammatone filterbank (Patterson et al., 1987) and a SNR is computed in each
274 frequency band.

275 In model 1, the advantages due to better-ear listening and binaural un-
276 masking are computed in parallel. In each time frame and frequency band, the
277 best SNR across the two ears is selected to obtain the better-ear SNR. The time
278 frame duration used for the better-ear SNR computation is 24 ms. A 20 dB ceil-
279 ing , which corresponds to the maximum SNR allowed in each frequency band
280 and time frame, is introduced in order to prevent the SNR to tend to infinity
281 in the temporal gaps of the masker. In each frequency band and time frame
282 the binaural unmasking advantage is computed using an equation proposed by
283 Culling et al. (2005) to estimate the binaural masking level difference (BMLD).
284 This value is computed using a time frame duration of 300 ms. The binaural
285 unmasking advantage and the better-ear SNR are then integrated across fre-
286 quency using the SII weightings (ANSI S3.5, 1997) and averaged across time
287 frames. The two values are added to obtain an effective SNR. This model is
288 able to accurately predict spatial release from masking (SRM) and dip listening
289 for speech presented against stationary and modulated noise maskers (Vicente
290 & Lavandier, 2020).

291 Model 2 was proposed by Prud'homme et al. (2022) and incorporates har-
292 monic cancellation. The masker is segmented into 300-ms time frames and the
293 mean F0 over the time frame is computed using PRAAT PSOLA (Boersma &
294 Weenink, 2018). A comb filter tuned to the masker F0 is applied to both target
295 and masker signals in order to simulate the mechanism of harmonic cancellation.
296 The harmonic cancellation component of the model adopted the same parame-
297 ters suggested by Prud'homme et al. (2020, 2022): a jitter in the estimation of
298 the F0 ($0.25F_0$), the width of the notches of the comb filter ($0.6F_0$), an SNR
299 ceiling (40 dB), and a frequency limit up to which harmonic cancellation is ap-
300 plied (5000 Hz). The better-ear SNR is computed (as in model 1) for both the
301 comb-filtered and unfiltered signals. The best better-ear SNR between the two
302 is chosen (i.e., harmonic cancellation is only applied if it provides an advantage).
303 Another condition for applying harmonic cancellation is that the masker signal
304 is voiced at least 50 % of the time in the considered time frame. If this is not
305 the case, the better-ear SNR and binaural unmasking advantage are computed

306 as in model 1. In model 2, harmonic cancellation and binaural unmasking are
307 mutually exclusive (if harmonic cancellation is active, only the better-ear SNR is
308 computed). The effective SNR is obtained by integrating the binaural unmask-
309 ing advantage and/or the better-ear SNR across frequency bands using the SII
310 weightings, and averaging them across time frames. As per Prud’homme et al.
311 (2020, 2022), to obtain the model predictions the model is run 800 times for
312 each condition using a different realization of the stimuli and a different value
313 of the jitter (drawn from a normal distribution with a standard deviation corre-
314 sponding to $0.25F_0$). The random jitter leads to a different prediction on each
315 of these “trials”, which are then averaged to obtain the final prediction. Model
316 2 proved useful to predict the effects of spatial separation, intonation and am-
317 plitude modulation for speech intelligibility against buzz maskers (Prud’homme
318 et al., 2022).

319 Both models are able to predict relative differences between SRTs in differ-
320 ent conditions, but cannot provide an absolute prediction of intelligibility. To
321 do so, a reference needs to be chosen, which is typically the average SRT across
322 conditions (Lavandier et al., 2012). The models were applied to the stimuli of
323 the three experiments of the present study. The reference chosen to fit the pre-
324 dictions to the data was the average SRT across conditions for each experiment.

325 *4.3. Results*

326 The gray lines in Figure 2 (top panel) present the model predictions for the
327 main experiment using model 1. This model does not predict any SRT difference
328 between SSN, monotonized buzzes and intonated buzzes. It underestimates the
329 difference in SRTs between SSN and the amplitude-modulated maskers (i.e., the
330 effect of dip listening). The model accurately predicts the SRM for SSN and
331 amplitude-modulated maskers, but it overestimates the SRM for buzz maskers.
332 The black lines in Figure 2 (top panel) present the model predictions using
333 model 2. This model accurately predicts the SRT differences between intonated
334 and monotonized buzzes. It accurately predicts the SRT differences between
335 buzzes and amplitude-modulated maskers. However, it underestimates the SRT

336 differences between SSN and the other masker types. The predicted SRM is
337 similar to that observed in the data, although it is slightly overestimated for
338 amplitude-modulated maskers.

339 Figure 2 (bottom panels) also presents the model predictions for the control
340 experiments using model 1 (gray lines) and model 2 (black lines). Both models
341 produce very similar predictions across conditions (less than 1 dB of variation),
342 consistent with the lack of any significant differences in the behavioral data.

343 5. Discussion

344 5.1. Harmonicity

345 In the main experiment, SRTs were highest for the SSN masker, which had
346 neither harmonicity nor the slow amplitude modulations that might support
347 temporal dip listening. SRTs were 5.9 and 11.5 dB higher for SSN compared
348 to intonated and monotonized buzzes, respectively, in the co-located condition.
349 This result is consistent with those of Steinmetzger & Rosen (2015), who mea-
350 sured SRTs for speech against SSN and intonated harmonic complexes and found
351 that SRTs were higher for SSN. The overall difference suggests that there is a
352 harmonicity-based benefit that could be due to harmonic cancellation and/or
353 spectral glimpsing. Another possibility is that periodic sounds cause less “mod-
354 ulation masking” than aperiodic sounds (Stone et al., 2011; Steinmetzger et al.,
355 2019).

356 Model 1, which should capture spectral glimpsing (by computing SNR within
357 frequency bands), predicts no differences between SSN, monotonized buzzes and
358 intonated buzzes, suggesting that differences in spectral glimpsing are negligible
359 in terms of SNR. Model 2, which incorporates harmonic cancellation, is able to
360 predict the difference between monotonized and intonated buzzes, consistent
361 with the results of Prud’homme et al. (2022). However, model 2 underestimates
362 the difference between SSN and buzzes by 4.6 dB. The only previous study that
363 tried to predict the SRT difference between noise and intonated buzz, using a
364 modulation-based model, also underestimated this effect by about 5 dB (Stein-

365 metzger et al., 2019). It seems that a combination of harmonic cancellation and
366 a mechanism that is sensitive to modulation masking may be required to fully
367 predict the benefits of harmonicity for simple maskers.

368 The finding of higher SRTs for intonated compared to monotonized buzzes
369 is consistent with a number of previous studies (Deroche & Culling, 2011; Green
370 & Rosen, 2013; Leclère et al., 2017). As proposed by Leclère et al. (2017), the
371 difference between intonated and monotonized maskers could be due to limita-
372 tions in the harmonic cancellation mechanism. Specifically, it may be that it is
373 difficult to “follow” the F0 contour when it varies over time. Another possibil-
374 ity is that the concomitant amplitude fluctuations that accompany intonation
375 increase modulation masking. It is also interesting to note that in the separated
376 conditions, the SRT difference between SSN and intonated buzz was greatly
377 reduced, to the point that there was no significant advantage of harmonicity.
378 It is possible that spatial separation provided enough masking release so that
379 there was no further benefit to be gained from the weak harmonicity cue in the
380 intonated condition.

381 No significant difference in SRT could be found between the speech and
382 monotonized speech maskers in any of the experiments. Given the robust effects
383 of intonation observed for buzzes, one might have expected a similar effect to
384 be observed for speech (i.e., lower SRTs for a monotonized speech masker than
385 for a naturally intonated speech masker). If the effect of intonation observed
386 for buzzes is due to harmonic cancellation, it apparently does not apply to
387 speech maskers. Another interesting result was that model 2 did not *predict*
388 an advantage for monotonized speech compared to naturally intonated speech,
389 contrary to the advantage it predicted for monotonized buzzes over intonated
390 buzzes. There was also no significant difference between the SRTs for speech and
391 vocoded speech. Given that the main difference between these two maskers is
392 harmonicity, this result provides a further indication that harmonicity does not
393 strongly affect the EM present in speech-on-speech situations. This result can
394 be compared to that of Rosen et al. (2013), who made a similar comparison but
395 did not attempt to reduce IM. In their one- and two-talker masker conditions,

396 performance was consistently poorer for speech maskers than for vocoded-speech
397 maskers (by about 3-7 dB). The fact that we estimated no differences in EM
398 bolsters the conclusion that their effects were related to IM.

399 So why might harmonic cancellation be so ineffective at reducing EM for
400 speech maskers? Our results point to the presence of unvoiced parts and/or
401 amplitude modulation in speech. The modeling results of Prud'homme et al.
402 (2022) with amplitude-modulated buzz maskers suggest that harmonic cancel-
403 lation can still be useful on amplitude-modulated signals. Thus, the most par-
404 simonious explanation for the lack of an effect of harmonic cancellation with
405 speech maskers may be the presence of unvoiced segments. Or perhaps it is a
406 combination of these effects. For example, if the target is primarily understood
407 based on information available in the temporal dips of a modulated masker,
408 and for speech maskers these low energy dips tend to be unvoiced rather than
409 voiced, this could explain why harmonic cancellation has little effect.

410 One concern with the absence of significant differences in SRT between the
411 different amplitude-modulated maskers in the main experiment was that the
412 SRTs were very low, and could have been affected by a floor effect. However,
413 the results from the two control experiments lessen this concern: measured SRTs
414 were substantially higher and the SRT differences between the four amplitude-
415 modulated masker conditions were still not significant. In control experiment 2,
416 our hypothesis was that if harmonic cancellation was at play, it would operate
417 most effectively for the monotonized speech masker with a single F0 as it should
418 be the easiest to cancel. The results do not confirm this hypothesis: the SRT
419 for this monotonized speech masker was not significantly different from that of
420 any other masker type.

421 Like in the main experiment, the control experiments revealed no significant
422 differences between SRTs for vocoded speech (with no harmonic structure) and
423 natural or monotonized speech. This suggests once again that harmonicity in
424 the masker did not play an important role here, and the model predictions are
425 consistent with that hypothesis. Of course, given the low number of participants
426 in these experiments, we cannot provide definitive evidence for the lack of an

427 effect of masker type. However, we assume that if there is any effect, it is very
428 small.

429 Our conclusion is that harmonicity-based effects on EM may be negligible
430 for speech-on-speech situations, and that examples of harmonicity-based release
431 from masking reported in the literature for speech maskers (Brokx & Noot-
432 boom, 1982; Deroche & Gracco, 2019; Popham et al., 2018) likely reflect a
433 release from IM. From a modeling perspective, this suggests that predictions of
434 EM in speech-on-speech situations might not need to take into account the ef-
435 fects of harmonicity, at least as a first approximation, even if it has been shown
436 to be important for masking caused by harmonic complexes (Prud’homme et al.,
437 2022). In other words, for predicting the EM present in cocktail party situations,
438 the “modulated-noise” models such as those proposed by Vicente & Lavandier
439 (2020) or Beutelmann et al. (2010) are likely to be sufficient.

440 5.2. *Spatial separation*

441 As expected, we observed a main effect of spatial separation between target
442 and masker, and the general pattern of SRM across masker types was consistent
443 with previous studies (Beutelmann & Brand, 2006; Jelfs et al., 2011; Leclère
444 et al., 2017; Culling & Lavandier, 2021).

445 The SRM predicted by model 1 is equivalent for SSN and buzzes, because this
446 model cannot distinguish buzzes from noises apart from their minor differences
447 in long-term spectrum (Fig. 1). This predicted SRM is similar to that observed
448 in the data for SSN, but larger than that observed for buzzes. For amplitude-
449 modulated maskers, the predicted SRM is in good agreement with the observed
450 SRM. The reason that the model predicts less SRM for amplitude-modulated
451 maskers appears to be due to the SNR ceiling set at 20 dB; as the dip listening
452 already provides a large SNR advantage, there is limited headroom for further
453 advantages.

454 Model 2 predicts the SRM reasonably well. Importantly, the success of this
455 model relies on the fact that harmonic cancellation and binaural unmasking are
456 mutually exclusive (Prud’homme et al., 2022). Without this assumption, the

457 model would overpredict SRM for harmonic maskers. The rationale behind this
458 assumption is that binaural unmasking is a release from simultaneous masking
459 and relies on the spectrotemporal overlap of the target and masker signals.
460 Thus, binaural unmasking should be larger for noise than for spectrally sparse
461 maskers like buzzes. Additional experimental work focused specifically on the
462 interaction between binaural unmasking and harmonic cancellation would be
463 needed for a more complete picture.

464 *5.3. Amplitude modulation*

465 All amplitude-modulated maskers (speech, monotonized speech, vocoded
466 speech, reversed speech) produced lower SRTs than the other masker types.
467 The difference in SRTs between SSN and vocoded speech was almost 15 dB in
468 the present study, which is larger than previously reported using similar stimuli
469 (e.g. Beutelmann et al., 2010; Collin & Lavandier, 2013). It may simply be that
470 our vocoded speech masker contained more or larger temporal dips than the
471 maskers used in previous studies. There is also some evidence that the benefit
472 derived from amplitude fluctuations in the masker depends on the nature of
473 the target speech materials (Schoof & Rosen, 2015; Best et al., 2019), and thus
474 the discrepancy may be explained by the use of matrix sentences (which are
475 potentially easier to guess than open-set materials) in our experiment. Indeed
476 model 1 has been shown to accurately predict dip listening benefits measured
477 with open-set sentences (Vicente et al., 2020), but it underpredicts the benefit
478 observed here with matrix sentences.

479 Leclère et al. (2017) found that masking release due to amplitude mod-
480 ulation was small for buzzes, to the point that there was no advantage for
481 monotonized buzzes with amplitude modulations compared to stationary mono-
482 tonized buzzes. Steinmetzger & Rosen (2015) also found that the advantage due
483 to amplitude modulation was smaller than the advantage due to periodicity. The
484 present results provide a slightly different picture for speech maskers, in that
485 evidence was found for temporal dip listening but no strong evidence for release
486 due to harmonicity. It is possible that the relative contribution of these two

487 mechanisms depends on the specific masker signal characteristics. For example,
488 harmonic cancellation may dominate for buzzes that have a rich harmonic struc-
489 ture, whereas dip listening may dominate for partially harmonic/voiced stimuli
490 such as speech.

491 **6. Conclusion**

492 SRTs were measured for maskers ranging from noise to speech in order to
493 better understand harmonicity-based contributions to EM in speech-on-speech
494 situations. The different masker types provided a comparison between maskers
495 with and without harmonic structure, amplitude modulation and variations in
496 F0 over time.

497 SRTs measured for unmodulated maskers (SSN and buzzes) suggest an ad-
498 vantage due to harmonicity in the masker that is impaired by intonation. Such
499 conditions continue to provide the most compelling case for a harmonic can-
500 cellation mechanism. On the other hand, the results for various amplitude-
501 modulated “speech-like” maskers suggest a very limited role for harmonic can-
502 cellation in reducing EM in mixtures of talkers. This conclusion is further sup-
503 ported by the predictions of models with and without a harmonic cancellation
504 component. Overall, we suggest that models validated for amplitude-modulated
505 noise maskers can be used as a first approximation for predicting EM in cocktail
506 party scenarios.

507 **7. Acknowledgments**

508 This work was performed within the LabEx CeLyA (Grant No. ANR-10-
509 LABX-0060) and funded by the “Fondation Pour l’Audition” (Speech2Ears
510 grant). V.B. was supported, in part, by National Institutes of Health-National
511 Institute on Deafness and Other Communication Disorders (NIH-NIDCD) Award
512 No. DC015760.

513 **References**

514 **References**

515 ANSI S3.5 (1997). Methods for Calculation of the Speech Intelligibility Index.
516 *American National Standards Institute, New York, .*

517 Best, V., Roverud, E., Baltzell, L., Rennies, J., & Lavandier, M. (2019). The
518 importance of a broad bandwidth for understanding “glimpsed” speech. *The*
519 *Journal of the Acoustical Society of America*, *146*, 3215. doi:10.1121/1.
520 5131651.

521 Beutelmann, R., & Brand, T. (2006). Prediction of speech intelligibility in
522 spatial noise and reverberation for normal-hearing and hearing-impaired lis-
523 teners. *The Journal of the Acoustical Society of America*, *120*, 331–342.
524 doi:10.1121/1.2202888.

525 Beutelmann, R., Brand, T., & Kollmeier, B. (2010). Revision, extension, and
526 evaluation of a binaural speech intelligibility model. *The Journal of the Acous-*
527 *tical Society of America*, *127*, 2479–2497. doi:10.1121/1.3295575.

528 Biberger, T., & Ewert, S. D. (2019). The effect of room acoustical parameters on
529 speech reception thresholds and spatial release from masking. *The Journal of*
530 *the Acoustical Society of America*, *146*, 2188–2200. doi:10.1121/1.5126694.

531 Boersma, P., & Weenink, D. (2018). Praat: Doing phonetics by com-
532 puter [Computer program]. Version 6.0.42, retrieved 15 August 2018 from
533 <http://www.praat.org/>.

534 Brokx, J. P., & Nootboom, S. G. (1982). Intonation and the perceptual sepa-
535 ration of simultaneous voices. *Journal of Phonetics*, *10*, 23–36.

536 Brungart, D. S., Simpson, B. D., Ericson, M. A., & Scott, K. R. (2001). Informa-
537 tional and energetic masking effects in the perception of multiple simultaneous
538 talkers. *The Journal of the Acoustical Society of America*, *110*, 2527–2538.
539 doi:10.1121/1.1408946.

- 540 Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with
541 One and with Two Ears. *The Journal of the Acoustical Society of America*,
542 25, 975–979. doi:10.1121/1.1907229.
- 543 Collin, B., & Lavandier, M. (2013). Binaural speech intelligibility in rooms
544 with variations in spatial location of sources and modulation depth of noise
545 interferers. *The Journal of the Acoustical Society of America*, 134, 1146–1159.
546 doi:10.1121/1.4812248.
- 547 Culling, J. F. (2016). Speech intelligibility in virtual restaurants. *The Journal of*
548 *the Acoustical Society of America*, 140, 2418–2426. doi:10.1121/1.4964401.
- 549 Culling, J. F., Hawley, M. L., & Litovsky, R. Y. (2005). Erratum: The role
550 head-induced interaural time and level differences in the speech reception
551 threshold for multiple interfering sound sources [J. Acoust. Soc. Am. 116,
552 1057 (2004)]. *The Journal of the Acoustical Society of America*, 118, 552–
553 552. doi:10.1121/1.1925967.
- 554 Culling, J. F., & Lavandier, M. (2021). Binaural unmasking and spatial release
555 from masking. In R. Y. Litovsky, M. J. Goupell, A. N. Popper, & R. R.
556 Fay (Eds.), *Binaural Hearing* (pp. 209–241). Springer Nature Switzerland
557 volume 73 of *Springer Handbook of Auditory Research*.
- 558 Darwin, C. J., Brungart, D. S., & Simpson, B. D. (2003). Effects of funda-
559 mental frequency and vocal-tract length changes on attention to one of two
560 simultaneous talkers. *The Journal of the Acoustical Society of America*, 114,
561 2913. doi:10.1121/1.1616924.
- 562 de Cheveigné, A. (1993). Separation of concurrent harmonic sounds: Fundamen-
563 tal frequency estimation and a time-domain cancellation model of auditory
564 processing. *The Journal of the Acoustical Society of America*, 93, 3271–3290.
565 doi:10.1121/1.405712.
- 566 de Cheveigné, A. (2021). Harmonic Cancellation—A Fundamental of Au-

- 567 ditory Scene Analysis. *Trends in Hearing*, 25, 2331–2165. doi:10.1177/
568 23312165211041422.
- 569 Deroche, M. L. D., & Culling, J. F. (2011). Voice segregation by difference in
570 fundamental frequency: Evidence for harmonic cancellation. *The Journal of*
571 *the Acoustical Society of America*, 130, 2855–2865. doi:10.1121/1.3643812.
- 572 Deroche, M. L. D., & Culling, J. F. (2013). Voice segregation by difference in
573 fundamental frequency: Effect of masker type. *The Journal of the Acoustical*
574 *Society of America*, 134, EL465–EL470. doi:10.1121/1.4826152.
- 575 Deroche, M. L. D., & Gracco, V. L. (2019). Segregation of voices with single
576 or double fundamental frequencies. *The Journal of the Acoustical Society of*
577 *America*, 145, 847–857. doi:10.1121/1.5090107.
- 578 Freyman, R. L., Balakrishnan, U., & Helfer, K. S. (2004). Effect of number
579 of masking talkers and auditory priming on informational masking in speech
580 recognition. *The Journal of the Acoustical Society of America*, 115, 2246–
581 2256. doi:10.1121/1.1689343.
- 582 Gardner, W. G., & Martin, K. D. (1995). HRTF measurements of a KEMAR.
583 *The Journal of the Acoustical Society of America*, 97, 3907–3908. doi:10.
584 1121/1.412407.
- 585 Glasberg, B. R., & Moore, B. C. J. (1990). Derivation of auditory filter
586 shapes from notched-noise data. *Hearing Research*, 47, 103–138. doi:10.
587 1016/0378-5955(90)90170-T.
- 588 Green, T., & Rosen, S. (2013). Phase effects on the masking of speech by
589 harmonic complexes: Variations with level. *The Journal of the Acoustical*
590 *Society of America*, 134, 2876–2883. doi:10.1121/1.4820899.
- 591 Jelfs, S., Culling, J. F., & Lavandier, M. (2011). Revision and validation of
592 a binaural model for speech intelligibility in noise. *Hearing Research*, 275,
593 96–104. doi:10.1016/j.heares.2010.12.005.

- 594 Kidd, G., Best, V., & Mason, C. R. (2008). Listening to every other word:
595 Examining the strength of linkage variables in forming streams of speech.
596 *The Journal of the Acoustical Society of America*, *124*, 3793–3802. doi:10.
597 1121/1.2998980.
- 598 Kidd, G., Mason, C. R., Swaminathan, J., Roverud, E., Clayton, K. K., &
599 Best, V. (2016). Determining the energetic and informational components of
600 speech-on-speech masking. *The Journal of the Acoustical Society of America*,
601 *140*, 132–144. doi:10.1121/1.4954748.
- 602 Lavandier, M., & Best, V. (2020). Modeling Binaural Speech Understand-
603 ing in Complex Situations. In J. Blauert, & J. Braasch (Eds.), *The Tech-*
604 *nology of Binaural Understanding* Modern Acoustics and Signal Process-
605 ing (pp. 547–578). Cham: Springer International Publishing. doi:10.1007/
606 978-3-030-00386-9_19.
- 607 Lavandier, M., Jelfs, S., Culling, J. F., Watkins, A. J., Raimond, A. P., & Makin,
608 S. J. (2012). Binaural prediction of speech intelligibility in reverberant rooms
609 with multiple noise sources. *The Journal of the Acoustical Society of America*,
610 *131*, 218–231. doi:10.1121/1.3662075.
- 611 Lavandier, M., Mason, C. R., Baltzell, L. S., & Best, V. (2021). Individual
612 differences in speech intelligibility at a cocktail party: A modeling perspective.
613 *The Journal of the Acoustical Society of America*, *150*, 1076–1087. doi:10.
614 1121/10.0005851.
- 615 Leclère, T., Lavandier, M., & Deroche, M. L. (2017). The intelligibility of speech
616 in a harmonic masker varying in fundamental frequency contour, broadband
617 temporal envelope, and spatial location. *Hearing Research*, *350*, 1–10. doi:10.
618 1016/j.heares.2017.03.012.
- 619 Patterson, R., Nimmo-Smith, I., Holdsworth, J., & Rice, P. (1987). An efficient
620 auditory filterbank based on the gammatone function. In *Presented to the*
621 *Institute of Acoustics Speech Group on Auditory Modelling at the Royal Signal*
622 *Research Establishment*. (p. 34).

- 623 Popham, S., Boebinger, D., Ellis, D. P. W., Kawahara, H., & McDermott,
624 J. H. (2018). Inharmonic speech reveals the role of harmonicity in the
625 cocktail party problem. *Nature Communications*, *9*, 2122. doi:10.1038/
626 s41467-018-04551-8.
- 627 Prud'homme, L., Lavandier, M., & Best, V. (2020). A harmonic-cancellation-
628 based model to predict speech intelligibility against a harmonic masker. *The*
629 *Journal of the Acoustical Society of America*, *148*, 3246–3254. doi:10.1121/
630 10.0002492.
- 631 Prud'homme, L., Lavandier, M., & Best, V. (2022). A dynamic binaural
632 harmonic-cancellation model to predict speech intelligibility against a har-
633 monic masker varying in intonation, temporal envelope, and location. *Sub-*
634 *mitted to this special issue, under review*, .
- 635 Rosen, S., Souza, P., Ekelund, C., & Majeed, A. A. (2013). Listening to speech
636 in a background of other talkers: Effects of talker number and noise vocoding.
637 *The Journal of the Acoustical Society of America*, *133*, 2431–2443. doi:10.
638 1121/1.4794379.
- 639 Rothauser, E. H., Chapman, W. D., Guttman, N., Nordby, K. S., Silbiger, H. R.,
640 Urbanek, G. E., & Weinstock, M. (1969). IEEE Recommended Practice for
641 Speech Quality Measurements. *IEEE Transactions on Audio and Electroac-*
642 *oustics*, *17*, 225–246. doi:10.1109/TAU.1969.1162058.
- 643 Schoof, T., & Rosen, S. (2015). High sentence predictability increases the fluc-
644 tuating masker benefit. *The Journal of the Acoustical Society of America*,
645 *138*, EL181–EL186. doi:10.1121/1.4929627.
- 646 Steinmetzger, K., & Rosen, S. (2015). The role of periodicity in perceiving
647 speech in quiet and in background noise. *The Journal of the Acoustical Society*
648 *of America*, *138*, 3586–3599. doi:10.1121/1.4936945.
- 649 Steinmetzger, K., Zaar, J., Relañó-Iborra, H., Rosen, S., & Dau, T. (2019).
650 Predicting the effects of periodicity on the intelligibility of masked speech:

- 651 An evaluation of different modelling approaches and their limitations. *The*
652 *Journal of the Acoustical Society of America*, *146*, 2562–2576. doi:10.1121/
653 1.5129050.
- 654 Stone, M. A., Füllgrabe, C., Mackinnon, R. C., & Moore, B. C. J. (2011).
655 The importance for speech intelligibility of random fluctuations in “steady”
656 background noise. *The Journal of the Acoustical Society of America*, *130*,
657 2874–2881. doi:10.1121/1.3641371.
- 658 Vicente, T., & Lavandier, M. (2020). Further validation of a binaural model
659 predicting speech intelligibility against envelope-modulated noises. *Hearing*
660 *Research*, *390*, 107937. doi:10.1016/j.heares.2020.107937.
- 661 Vicente, T., Lavandier, M., & Buchholz, J. M. (2020). A binaural model im-
662 plementing an internal noise to predict the effect of hearing impairment on
663 speech intelligibility in non-stationary noises. *The Journal of the Acoustical*
664 *Society of America*, *148*, 3305–3317. doi:10.1121/10.0002660.
- 665 Wasiuk, P. A., Lavandier, M., Buss, E., Oleson, J., & Calandruccio, L. (2020).
666 The effect of fundamental frequency contour similarity on multi-talker listen-
667 ing in older and younger adults. *The Journal of the Acoustical Society of*
668 *America*, *148*, 3527–3543. doi:10.1121/10.0002661.
- 669 Wichmann, F. A., & Hill, N. J. (2001). The psychometric function: I. Fitting,
670 sampling, and goodness of fit. *Perception & Psychophysics*, *63*, 1293–1313.
671 doi:10.3758/BF03194544.