

1 **(Febr 1, 2018) ACCEPTED FOR PUBLICATION IN:**
2 **JOURNAL OF THE ACOUSTICAL SOCIETY OF AMERICA**
3 **EXPRESS LETTER**

4
5 **Talkers produce more pronounced amplitude modulations when**
6 **speaking in noise**

7
8 **Hans Rutger Bosker¹**

9 *Max Planck Institute for Psycholinguistics, PO Box 310, 6500 AH, Nijmegen, The Netherlands*
10 *Donders Institute for Brain, Cognition and Behaviour, Radboud University, Nijmegen, the Netherlands*
11 HansRutger.Bosker@mpi.nl

12
13 **Martin Cooke**

14 *Language and Speech Laboratory, Universidad del País Vasco, Vitoria, 01006, Spain*
15 *Ikerbasque (Basque Science Foundation), Bilbao, Spain*
16 m.cooke@ikerbasque.org

17
18 **Word count: 4266**

19 **Table count: 1**

20 **Figure count: 2**

21

¹ Corresponding author. Tel.: +31 (0)24 3521 373.
E-mail address: HansRutger.Bosker@mpi.nl

22

ABSTRACT

23 Speakers adjust their voice when talking in noise (known as Lombard speech), facilitating
24 speech comprehension. Recent neurobiological models of speech perception emphasize the role
25 of amplitude modulations in speech-in-noise comprehension, helping neural oscillators to ‘track’
26 the attended speech. This study tested whether talkers produce more pronounced amplitude
27 modulations in noise. Across four different corpora, modulation spectra showed greater power in
28 amplitude modulations below 4 Hz in Lombard speech compared to matching plain speech. This
29 suggests that noise-induced speech contains more pronounced amplitude modulations, potentially
30 helping the listening brain to entrain to the attended talker, aiding comprehension.

31

32 *Keywords:* speech in noise, Lombard speech, temporal envelope, amplitude modulations,
33 modulation spectrum.

34

35

INTRODUCTION

36 Speakers typically adjust their voice when talking in noisy conditions. Speech produced in
37 noise generally exhibits increased intensity, slower speech rate, raised F0 and flatter spectral tilt
38 (for an overview, see Cooke, King, Garnier, & Aubanel, 2014). These and other modifications
39 result in what is collectively known as Lombard speech (Lombard, 1911). The scientific
40 importance of this form of speech stems in large part from the discovery by Dreher and O'Neill
41 (1957) that, after discounting intensity increases, Lombard speech is more intelligible than
42 unmodified speech when presented in noise; a finding that has been confirmed in a number of
43 subsequent studies (e.g., Pittman & Wiley, 2001; Summers, Pisoni, Bernacki, Pedlow, & Stokes,
44 1988). However, the basis for intelligibility gains is not fully understood. One aspect of noise-
45 induced speech that has received little attention concerns how talkers adjust the temporal
46 modulations of their speech when conversing in noise. The present study examines the temporal
47 modulations in Lombard speech and plain speech (speech produced in quiet) and demonstrates
48 that amplitude modulations are enhanced in the temporal envelope of Lombard speech compared
49 to matching plain speech.

50 Speech is an inherently rhythmic signal in that it contains strong amplitude modulations,
51 particularly in the 1-15 Hz range (Ding et al., 2017; Varnet, Ortiz-Barajas, Erra, Gervain, &
52 Lorenzi, 2017). These amplitude modulations, evident in the temporal envelope of speech,
53 greatly contribute to speech intelligibility (Drullman, Festen, & Plomp, 1994; Ghitza, 2012;
54 Shannon, Zeng, Kamath, Wygonski, & Ekelid, 1995; Smith, Delgutte, & Oxenham, 2002).
55 Speech with more pronounced amplitude modulations is more intelligible in noise (Houtgast &
56 Steeneken, 1985; Koutsogiannaki & Stylianou, 2016; Steeneken & Houtgast, 1980). Also,

57 speakers who are intrinsically more intelligible than others show more pronounced low-
58 frequency modulations in the temporal envelope (Bradlow, Torretta, & Pisoni, 1996).
59 Modulations as low as 2 Hz have been shown to be essential for phoneme identification
60 (Drullman et al., 1994). In fact, removing amplitude modulations, occurring at a syllabic rate,
61 from speech impairs its intelligibility to a large degree (Ghitza, 2012).

62 Recent neurobiological models of speech perception (Ghitza, 2011; Giraud & Poeppel, 2012;
63 Peelle & Davis, 2012) propose that enhanced temporal modulations facilitate speech perception
64 because speech-envelope information evokes marked “envelope-following” neural responses in
65 the auditory cortex, known as *neural entrainment*. Endogenous neural oscillators in the delta (1-4
66 Hz) and theta range (4-8 Hz) are thought to phase-lock (entrain) to the amplitude fluctuations in
67 the input signal (Doelling, Arnal, Ghitza, & Poeppel, 2014). Thus, neuronal excitability is
68 temporally aligned with the temporal structure of the attended spoken input, serving as a parsing
69 mechanism for the initial neural coding of the speech signal (Arnal, Giraud, & Poeppel, 2015;
70 Bosker, 2017; Bosker & Kösem, 2017; Kösem et al., 2017).

71 This neural entrainment to the temporal modulations in speech has been proposed to be one of
72 the mechanisms by which listeners are capable of understanding speech in challenging listening
73 conditions, such as in noise or with competing talkers. Brain oscillations during speech-in-noise
74 perception preferentially track attended relative to ignored speech streams, using particularly the
75 phase of low-frequency neural activity (1–8 Hz) (Ding & Simon, 2012; Kerlin, Shahin, & Miller,
76 2010). The intelligibility of an attended speech stream in noise can be predicted from the extent
77 to which cortical oscillators are aligned to the temporal envelope of the attended signal
78 (Golumbic et al., 2013; Golumbic, Poeppel, & Schroeder, 2012; Rimmele, Golumbic, Schröger,
79 & Poeppel, 2015). In fact, modulating listeners’ neural activity with transcranial stimulation with

80 speech-envelope-shaped currents has been argued to help speech-in-noise comprehension
81 (Riecke, Formisano, Sorger, Başkent, & Gaudrain, 2018).

82 Considering these neurobiological models and the reported beneficial effects of enhanced
83 amplitude modulations on perception, it may be hypothesized that speakers, in an attempt to aid
84 speech intelligibility, would also naturally produce more enhanced amplitude modulations when
85 talking in a noisy acoustic environment. This would allow greater opportunity for the listening
86 brain to entrain to the enhanced temporal envelope. There is some evidence for larger within-
87 syllable intensity changes (Garnier & Henrich, 2014) and greater overall RMS range (Folk &
88 Schiel, 2011) in Lombard speech (relative to speech-in-quiet) but none of these studies examined
89 the strength of amplitude modulations in the temporal envelope.

90 Krause and Braida (2004) investigated another kind of speech adjustment, namely *clear*
91 *speech*. In contrast to Lombard speech, clear speech is elicited in quiet environments by
92 explicitly asking speakers to speak more clearly (e.g., by imagining speaking to a hearing-
93 impaired person; Uchanski, 2008). Modulation spectra, showing the power of frequency
94 components in the temporal envelope, revealed stronger amplitude modulations below 4 Hz in
95 clear speech. This difference between clear and plain speech was most apparent in frequency
96 bands around 500, 1000, and 2000 Hz. However, the effect was only observed for 2 talkers (T3
97 and T5; 50 trials each) and was induced through instructions rather than by physically presented
98 adverse listening conditions (Krause & Braida, 2002).

99 Therefore, the present study examines the temporal modulations in Lombard speech and plain
100 speech, adapting the method from Krause and Braida (2004). Using modulation spectra, the
101 power of amplitude modulations in the temporal envelope of Lombard speech and plain speech is
102 compared in three modulation frequency bands: the delta range (1-4 Hz), the theta range (4-8

103 Hz), and the alpha range (8-15 Hz)². In neurobiological studies, neural entrainment is particularly
104 observed in the lower frequency range. Accordingly, we hypothesize that, across several speech
105 corpora, Lombard speech will have more pronounced amplitude modulations compared to plain
106 speech in the delta/theta range, as evidenced by greater power in the modulation spectrum.

107 **METHOD**

108 Four English speech corpora were analyzed (see Table 1), each including sentences produced
109 in quiet and the same sentences produced in noisy elicitation conditions. The corpora varied on
110 several dimensions: for instance, Corpus 1 and 3 used ‘normal’ sentences (i.e., meaningful
111 everyday sentences), while Corpus 2 used six-word matrix sentences (e.g., “lay green with A4
112 now” or “set white at B8 again”) and Corpus 4 used frame sentences (e.g., “Now we will say
113 CV(C)C again”). Corpora also varied in the noise conditions and loudness levels used to elicit
114 Lombard speech; for instance, some used speech-shaped noise (i.e., noise with speech-like
115 LTAS), others used noise modulated by a single talker or multiple talkers (e.g., ICRA noise from
116 Dreschler, Verschuure, Ludvigsen, & Westermann, 2001).

117

² Note that we use the terms delta, theta, and alpha to refer to specific modulation frequency bands (1-4 Hz; 4-8 Hz; 8-15 Hz, respectively) irrespective of whether these modulations occur in speech or in neural signals.

118

119 Table 1. Characteristics of the four speech corpora (M = male; F = female; BMN = 9-talker babble-modulated ICRA
120 noise from Dreschler et al. (2001); SSN = speech-shaped noise; SMN = speech-modulated noise).

	Talkers	Sentences	Noise	Source
Corpus 1	$N = 1$ (M)	‘normal’; $N = 25$	BMN; intense	Mayo, Aubanel, and Cooke (2012); doi: 10.7488/ds/138
Corpus 2	$N = 8$ (4 M/4 F)	matrix $N = 400$	SSN; 96 dB (L)	Lu and Cooke (2008)
Corpus 3 (Hurricane)	$N = 1$ (M)	‘normal’; $N = 720$	SMN; 84 dB (A)	Cooke et al. (2013); doi: 10.7488/ds/140
Corpus 4 (MRT)	$N = 1$ (M)	frame; $N = 300$	SMN; 84 dB (A)	Collected by Cassia Valentino-Botinhao; http://datashare.is.ed.ac.uk/handle/10283/347

121

122 Before analysis, any leading and trailing silences around the sentences were manually
123 removed. Two types of analysis were performed: a broadband analysis and a filterbank analysis.
124 Both the broadband analysis and the filterbank analysis involved calculating the modulation
125 spectrum of each sentence in each corpus using a method adapted from Krause and Braidia
126 (2004). This included normalizing the overall power of signals (root-mean-square; RMS),
127 matching the overall energy of plain and Lombard signals. Thus, any potential differences
128 between plain and Lombard speech cannot be attributed to differences in overall energy.

129 For the broadband analysis, each sentence was filtered by a sixth-order Butterworth band-pass
130 filter spanning the 250-4000 Hz range, followed by estimation of the envelope of the filter’s
131 output via the Hilbert transform. The envelope signal was then submitted to a Fast Fourier

132 Transform, resulting in the modulation spectrum of that particular speech fragment. Finally, for
133 statistical comparisons, the average power in three frequency bands was calculated: average
134 power in the 1-4 Hz range (delta), the 4-8 Hz range (theta), and the 8-15 Hz range (alpha). This
135 resulted in three different observations for each sentence, forming the dependent variables for the
136 statistical analyses reported below.

137 The filterbank analysis was performed to further investigate whether any potential difference
138 between the plain and Lombard speech could be attributed to particular frequency bands. The
139 filterbank analysis was identical to the broadband analysis, except that the speech signal was
140 filtered into five component signals, using a bank of fourth-order Butterworth filters with center
141 frequencies of 500 Hz (bandwidth: 125 Hz), 1000 Hz (250 Hz), 2000 Hz (500 Hz), 4000 Hz
142 (1000 Hz), and 8000 Hz (2000 Hz). The bandlimited output of this filterbank formed the input
143 for the subsequent calculation of the modulation spectrum separately for each frequency band
144 (using the procedure described above).

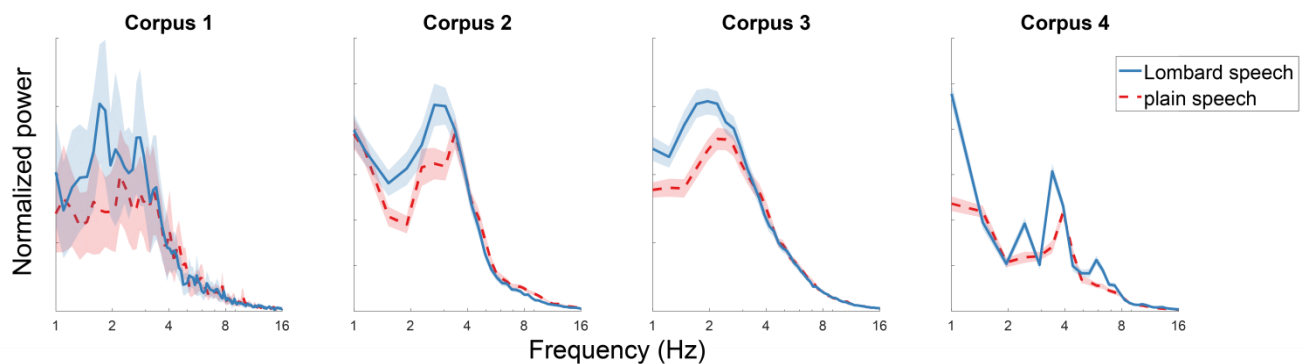
145 Since each sentence had its own unique duration, this procedure resulted in modulation spectra
146 with unique frequency resolutions. However, in order to visualize the *average* rhythmicity in
147 plain and Lombard speech across multiple sentences, identical frequency resolutions were
148 required. This was achieved by repeating the procedure above with the envelope signal zero-
149 padded to the next power of 2 higher than the length of the longest fragment of that particular
150 corpus. Note that this zero-padding was only performed for visualization purposes; statistical
151 analyses were performed on the data from the non-padded signals.

152 **RESULTS**

153 Figure 1 shows the average modulation spectra from the broadband analysis (250-4000 Hz)

154 for plain and Lombard speech, for each corpus.

155



156

157 Figure 1. (color online) Average modulation spectra, calculated from the broadband analysis (250-4000 Hz), of
158 Lombard speech (solid line) and matching plain speech (dashed line), for each corpus. Shaded areas indicate 95%
159 CIs.

160

161 The first thing to note is that the modulation spectra of Lombard speech resemble those of
162 plain speech in overall shape (e.g., number of peaks and troughs), indicating that the temporal
163 structure of matching Lombard and plain utterances is very much alike. Secondly, the peaks in
164 the modulation spectra of Lombard speech occur at slightly lower frequencies (i.e., shifted
165 leftward) than the peaks in the modulation spectra of plain speech. This observation is in line
166 with the fact that Lombard speech typically has a slower speech rate than plain speech, shifting
167 the temporal modulations down towards lower frequencies. In fact, the size of the shift observed
168 in the modulation spectra is in keeping with the average rate change in Lombard speech (e.g.,
169 7.6% in Corpus 2; Lu & Cooke, 2008).

170 Finally, there would seem to be consistently higher power in Lombard speech, across the four
171 corpora, in the lower frequency range between 1-4 Hz (delta). This difference was statistically
172 analyzed by means of Linear Mixed Models (LMMs; Baayen, 2008) as implemented in the lme4

173 library in R. Our three dependent variables, average power in delta (1-4 Hz), theta (4-8 Hz), and
174 alpha range (8-15 Hz), were entered into separate LMMs with identical structure. Condition
175 (categorical predictor, with the plain condition mapped onto the intercept), Corpus (categorical
176 predictor, with Corpus 1 mapped onto the intercept), and their interaction, were entered as
177 predictors, with Talker entered as random factor with by-talker random slopes for Condition
178 (Barr, Levy, Scheepers, & Tily, 2013). More complex models including a by-talker random slope
179 for Condition failed to converge. Statistical significance was assessed at the 0.05 significance
180 level by checking whether effects had absolute t -values exceeding 2 (Baayen, 2008).

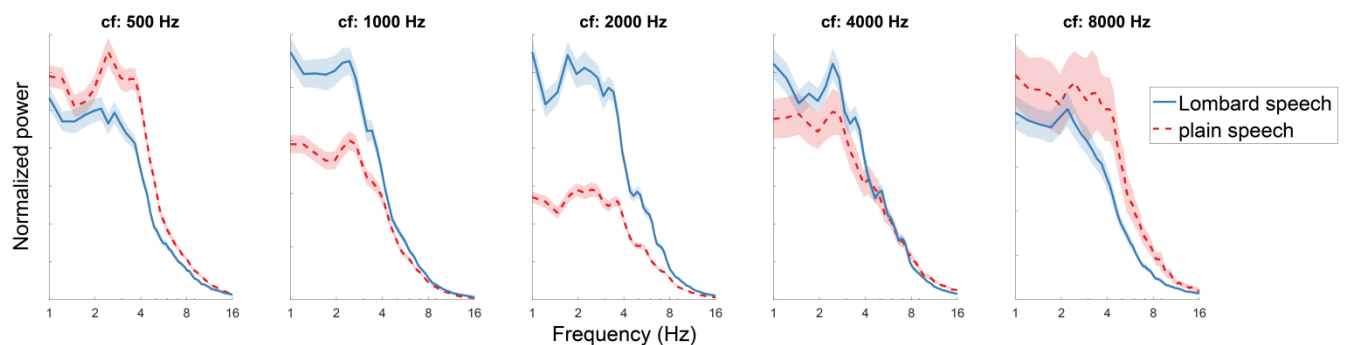
181 Only the model of average power *in the delta range* (marginal $R^2 = 0.185$; conditional $R^2 =$
182 0.408) revealed a significant effect: Lombard speech had a significantly higher average power in
183 the delta range compared to plain speech ($\beta = 0.008$, $SE = 0.002$, $t = 3.855$). The fact that no
184 significant interactions between Condition and Corpus were observed suggests that the effect of
185 Condition held equally across all four corpora. No significant effects or interactions were
186 observed in the other two models.

187 In order to explore whether this difference in overall power between Lombard and plain
188 speech in the delta range happened to be a by-product of the lower average speech rate in
189 Lombard speech, a follow-up analysis was performed. This analysis involved artificially
190 matching the (slower) speech rate of the Lombard sentences to the (faster) speech rate of the
191 plain sentences. We adopted the global duration modifications described in Cooke, Mayo, and
192 Villegas (2014), involving linear compression of the Lombard utterances using PSOLA in Praat.
193 The results of statistical analyses of the modulation spectra of plain and these duration-matched
194 Lombard sentences mirrored the results reported above. We found no significant effects in the
195 theta or alpha range, and only one significant effect of Condition in the delta range ($\beta = 0.011$,

196 $SE = 0.005$, $t = 2.230$; model's marginal $R^2 = 0.127$; conditional $R^2 = 0.226$), corroborating that
 197 duration-matched Lombard utterances had higher overall power in the delta range than plain
 198 utterances across the four corpora.

199 The filterbank analysis was designed to further investigate which frequency bands drive the
 200 difference in temporal modulations between plain and Lombard speech. Figure 2 shows the
 201 average modulation spectra for plain and Lombard speech from the filterbank analysis (with
 202 center frequencies at 500-8000 Hz), averaging over the four different corpora.

203



204

205 Figure 2. (color online) Average modulation spectra, calculated from the filterbank analysis (center frequencies at
 206 500-8000 Hz), of Lombard speech (solid line) and matching plain speech (dashed line), averaging over the different
 207 corpora (cf = center frequency). Shaded areas indicate 95% CIs.

208

209 Judging from Figure 2, higher overall power in Lombard speech in the delta range would seem
 210 to be primarily driven by the 1000 Hz and 2000 Hz bands. Differences between Lombard and
 211 plain speech in the delta range were statistically tested by means of another LMM. This LMM
 212 predicted average power in the delta range (1-4 Hz) with fixed effects of Condition (categorical
 213 predictor, with the plain condition mapped onto the intercept), Band (categorical predictor, with

214 the 2000 Hz frequency band mapped onto the intercept), and their interaction. Talker was entered
215 as random factor with by-talker random slopes for Condition and Band (marginal $R^2 = 0.686$;
216 conditional $R^2 = 0.773$). More complex models including the predictor Corpus failed to
217 converge.

218 A simple effect of Condition revealed significantly higher average power in the Lombard
219 condition (relative to plain) in the delta range in the 2000 Hz frequency band (being mapped onto
220 the intercept; $\beta = 0.003$, $SE = 0.0002$, $t = 14.850$). Interactions between Condition and Band led
221 to two observations: first, the absence of an interaction between Condition and the 1000 Hz
222 frequency band showed that the effect of Condition was comparable in the 1000 Hz and 2000 Hz
223 frequency bands. Second, the effect of Condition was considerably smaller or even absent in the
224 500 Hz, 4000 Hz, and 8000 Hz frequency bands (500 Hz: $\beta = -0.004$, $SE = 0.0001$, $t = -22.815$;
225 4000 Hz: $\beta = -0.003$, $SE = 0.0002$, $t = -14.705$; 8000 Hz: $\beta = -0.003$, $SE = 0.0002$, $t = -16.868$).

226

DISCUSSION

227 This study compared the power of amplitude modulations in the temporal envelope of
228 Lombard speech (sentences produced in noise) and plain speech (same sentences but produced in
229 quiet). Speech from four different corpora (various speakers, sentences types, elicitation
230 methods; cf. Table 1) was analyzed by means of modulation spectra, revealing the power of
231 frequency components in the temporal envelopes. Across all four corpora, amplitude modulations
232 below 4 Hz were stronger in Lombard speech compared to matched plain speech (cf. similar
233 findings for clear speech; Krause & Braida, 2002, 2004). This difference was shown to be
234 independent of changes in speech rate and was concentrated in frequency bands around 1000 Hz
235 and 2000 Hz (cf. similar findings for clear speech; Krause & Braida, 2002, 2004).

236 The modulations most affected by Lombard speech (i.e., in delta range; 1-4 Hz) correspond to
237 the average syllable rates (e.g., Corpus 2; plain: 3.7 Hz; Lombard: 3.4 Hz). This suggests that the
238 difference in amplitude modulations in Lombard and plain speech may be driven by more
239 pronounced syllabic energy fluctuations in Lombard speech. The effect was concentrated in
240 higher frequency bands (around 1000 Hz and 2000 Hz), which is in line with studies on artificial
241 speech enhancement. For instance, Koutsogiannaki and Stylianou (2016) found that artificially
242 decreasing the modulation depth in lower frequency regions (200-600 Hz) and increasing the
243 modulation depth in higher frequencies (800-3000 Hz) enhanced speech-in-noise intelligibility.

244 These results suggest that speakers produce more pronounced amplitude modulations in noise
245 compared to in quiet. We interpret these findings in light of recent oscillations-based models of
246 speech perception (Ghitza, 2011; Giraud & Poeppel, 2012; Peelle & Davis, 2012), whereby
247 neural oscillations phase-lock (entrain) to amplitude fluctuations in speech. Greater amplitude
248 modulations in speech produced in noise presumably help the listening brain to entrain to the
249 attended talker, aligning neuronal excitability to the temporal structure of the attended signal,
250 facilitating speech-in-noise perception.

251 This idea is corroborated by neurobiological studies showing that removing amplitude
252 modulations, occurring at the syllabic rate, from speech reduces neural envelope-tracking
253 activity, impairing intelligibility (Doelling et al., 2014). Similarly, perception studies have shown
254 beneficial effects of temporal modulations on phoneme identification (Drullman et al., 1994) and
255 intelligibility (Ghitza, 2012). Future neuroimaging studies should investigate whether the
256 observed greater modulation depth in Lombard speech indeed facilitates cortical speech-tracking,
257 aiding speech-in-noise intelligibility.

258

ACKNOWLEDGEMENTS

259 The first author was supported by a Gravitation grant from the Dutch Government to the
260 Language in Interaction Consortium.

261

REFERENCES

262 Arnal, L. H., Giraud, A.-L., & Poeppel, D. (2015). A Neurophysiological Perspective on Speech Processing in "The
263 Neurobiology of Language". In G. Hickok & S. Small (Eds.), *Neurobiology of Language* (pp. 463-478).
264 San Diego: Academic Press.

265 Baayen, R. H. (2008). *Analyzing linguistic data: A practical introduction to statistics using R*. Cambridge:
266 Cambridge University Press.

267 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis
268 testing: Keep it maximal. *Journal of Memory and Language*, 68(3), 255-278.

269 Bosker, H. R. (2017). Accounting for rate-dependent category boundary shifts in speech perception. *Attention,*
270 *Perception & Psychophysics*, 79(1), 333-343. doi:10.3758/s13414-016-1206-4

271 Bosker, H. R., & Kösem, A. (2017). *An entrained rhythm's frequency, not phase, influences temporal sampling of*
272 *speech*. Paper presented at the Proceedings of Interspeech 2017, Stockholm.

273 Bradlow, A. R., Torretta, G. M., & Pisoni, D. B. (1996). Intelligibility of normal speech I: Global and fine-grained
274 acoustic-phonetic talker characteristics. *Speech Communication*, 20(3-4), 255-272.

275 Cooke, M., King, S., Garnier, M., & Aubanel, V. (2014). The listening talker: A review of human and algorithmic
276 context-induced modifications of speech. *Computer Speech & Language*, 28(2), 543-571.

277 Cooke, M., Mayo, C., Valentini-Botinhao, C., Stylianou, Y., Sauert, B., & Tang, Y. (2013). Evaluating the
278 intelligibility benefit of speech modifications in known noise conditions. *Speech Communication*, 55(4),
279 572-585.

280 Cooke, M., Mayo, C., & Villegas, J. (2014). The contribution of durational and spectral changes to the Lombard
281 speech intelligibility benefit. *The Journal of the Acoustical Society of America*, 135(2), 874-883.

282 Ding, N., Patel, A., Chen, L., Butler, H., Luo, C., & Poeppel, D. (2017). Temporal modulations in speech and music.

- 283 *Neuroscience and Biobehavioral Reviews, Online version*. doi:10.1016/j.neubiorev.2017.02.011
- 284 Ding, N., & Simon, J. Z. (2012). Emergence of neural encoding of auditory objects while listening to competing
285 speakers. *Proceedings of the National Academy of Sciences, 109*(29), 11854-11859.
- 286 Doelling, K. B., Arnal, L. H., Ghitza, O., & Poeppel, D. (2014). Acoustic landmarks drive delta–theta oscillations to
287 enable speech comprehension by facilitating perceptual parsing. *NeuroImage, 85*, 761-768.
- 288 Dreher, J. J., & O'Neill, J. (1957). Effects of ambient noise on speaker intelligibility for words and phrases. *The*
289 *Journal of the Acoustical Society of America, 29*(12), 1320-1323.
- 290 Dreschler, W. A., Verschuure, H., Ludvigsen, C., & Westermann, S. (2001). ICRA noises: Artificial noise signals
291 with speech-like spectral and temporal properties for hearing instrument assessment. *Audiology, 40*(3), 148-
292 157.
- 293 Drullman, R., Festen, J. M., & Plomp, R. (1994). Effect of reducing slow temporal modulations on speech
294 recognition.” *J. Acoust. Soc. Am, 95*(5), 2670-2680.
- 295 Folk, L., & Schiel, F. (2011). The Lombard Effect in spontaneous dialog speech. *Proceedings of Interspeech* (pp.
296 2701-2704). Florence.
- 297 Garnier, M., & Henrich, N. (2014). Speaking in noise: How does the Lombard effect improve acoustic contrasts
298 between speech and ambient noise? *Computer Speech & Language, 28*(2), 580-597.
- 299 Ghitza, O. (2011). Linking speech perception and neurophysiology: speech decoding guided by cascaded oscillators
300 locked to the input rhythm. *Frontiers in Psychology, 2*(130).
- 301 Ghitza, O. (2012). On the role of theta-driven syllabic parsing in decoding speech: intelligibility of speech with a
302 manipulated modulation spectrum. *Frontiers in Psychology, 3*, 238.
- 303 Giraud, A.-L., & Poeppel, D. (2012). Cortical oscillations and speech processing: emerging computational principles
304 and operations. *Nature Neuroscience, 15*(4), 511-517.
- 305 Golumbic, E. M. Z., Ding, N., Bickel, S., Lakatos, P., Schevon, C. A., McKhann, G. M., . . . Simon, J. Z. (2013).
306 Mechanisms underlying selective neuronal tracking of attended speech at a “cocktail party”. *Neuron, 77*(5),
307 980-991.
- 308 Golumbic, E. M. Z., Poeppel, D., & Schroeder, C. E. (2012). Temporal context in speech processing and attentional
309 stream selection: A behavioral and neural perspective. *Brain and Language, 122*(3), 151-161.

- 310 Houtgast, T., & Steeneken, H. J. (1985). A review of the MTF concept in room acoustics and its use for estimating
311 speech intelligibility in auditoria. *The Journal of the Acoustical Society of America*, 77(3), 1069-1077.
- 312 Kerlin, J. R., Shahin, A. J., & Miller, L. M. (2010). Attentional gain control of ongoing cortical speech
313 representations in a “cocktail party”. *The Journal of Neuroscience*, 30(2), 620-628.
- 314 Kösem, A., Bosker, H. R., Takashima, A., Jensen, O., Meyer, A., & Hagoort, P. (2017). Neural entrainment
315 determines the words we hear. *bioRxiv*. doi:10.1101/175000
- 316 Koutsogiannaki, M., & Stylianou, Y. (2016). Modulation Enhancement of Temporal Envelopes for Increasing
317 Speech Intelligibility in Noise. *Proceedings of Interspeech* (pp. 2508-2512).
- 318 Krause, J. C., & Braida, L. D. (2002). Investigating alternative forms of clear speech: The effects of speaking rate
319 and speaking mode on intelligibility. *The Journal of the Acoustical Society of America*, 112(5), 2165-2172.
- 320 Krause, J. C., & Braida, L. D. (2004). Acoustic properties of naturally produced clear speech at normal speaking
321 rates. *The Journal of the Acoustical Society of America*, 115(1), 362-378.
- 322 Lu, Y., & Cooke, M. (2008). Speech production modifications produced by competing talkers, babble, and
323 stationary noise. *The Journal of the Acoustical Society of America*, 124(5), 3261-3275.
- 324 Mayo, C., Aubanel, V., & Cooke, M. (2012). Effect of prosodic changes on speech intelligibility *Proceedings of*
325 *Interspeech* (pp. 1708-1711). Portland.
- 326 Peelle, J. E., & Davis, M. H. (2012). Neural oscillations carry speech rhythm through to comprehension. *Frontiers*
327 *in Psychology*, 3. doi:10.3389/fpsyg.2012.00320
- 328 Pittman, A. L., & Wiley, T. L. (2001). Recognition of speech produced in noise. *Journal of Speech, Language, and*
329 *Hearing Research*, 44(3), 487-496.
- 330 Riecke, L., Formisano, E., Sorger, B., Başkent, D., & Gaudrain, E. (2018). Neural entrainment to speech modulates
331 speech intelligibility. *Current Biology*. doi:10.1016/j.cub.2017.11.033
- 332 Rimmele, J. M., Golumbic, E. M. Z., Schröger, E., & Poeppel, D. (2015). The effects of selective attention and
333 speech acoustics on neural speech-tracking in a multi-talker scene. *Cortex*, 68, 144-154.
- 334 Shannon, R. V., Zeng, F.-G., Kamath, V., Wygonski, J., & Ekelid, M. (1995). Speech recognition with primarily
335 temporal cues. *Science*, 270(5234), 303.
- 336 Smith, Z. M., Delgutte, B., & Oxenham, A. J. (2002). Chimaeric sounds reveal dichotomies in auditory perception.

337 *Nature*, 416(6876), 87-90.

338 Steeneken, H. J., & Houtgast, T. (1980). A physical method for measuring speech-transmission quality. *The Journal*
339 *of the Acoustical Society of America*, 67(1), 318-326.

340 Summers, W. V., Pisoni, D. B., Bernacki, R. H., Pedlow, R. I., & Stokes, M. A. (1988). Effects of noise on speech
341 production: Acoustic and perceptual analyses. *The Journal of the Acoustical Society of America*, 84(3),
342 917-928.

343 Uchanski, R. M. (2008). Clear Speech. *The handbook of speech perception* (pp. 207-235): Blackwell Publishing Ltd.

344 Varnet, L., Ortiz-Barajas, M. C., Erra, R. G., Gervain, J., & Lorenzi, C. (2017). A cross-linguistic study of speech
345 modulation spectra. *The Journal of the Acoustical Society of America*, 142(4), 1976-1989.

346

347

348

LIST OF TABLES

349 Table 1. Characteristics of the four speech corpora (M = male; F = female; BMN = 9-talker babble-modulated ICRA
350 noise from Dreschler et al. (2001); SSN = speech-shaped noise; SMN = speech-modulated
351 noise)..... p. 7

352

353

LIST OF FIGURES

354 Figure 1. (color online) Average modulation spectra, calculated from the broadband analysis (250-4000 Hz), of
355 Lombard speech (solid line) and matching plain speech (dashed line), for each corpus. Shaded areas indicate 95%
356 CIs..... p. 8

357 Figure 2. (color online) Average modulation spectra, calculated from the filterbank analysis (center frequencies at
358 500-8000 Hz), of Lombard speech (solid line) and matching plain speech (dashed line), averaging over the different
359 corpora (cf = center frequency). Shaded areas indicate 95% CIs..... p. 9

360