

# Speech produced in noise: Relationship between listening difficulty and acoustic and durational parameters<sup>a)</sup>

Simone Graetzer,<sup>1,b)</sup> Pasquale Bottalico,<sup>2</sup> and Eric J. Hunter<sup>2</sup>

<sup>1</sup>Acoustics Research Unit, School of Architecture, University of Liverpool, Liverpool, England

<sup>2</sup>Voice Biomechanics and Acoustics Laboratory, Department of Communicative Sciences and Disorders, Michigan State University, East Lansing, Michigan 48824, USA

(Received 10 February 2016; revised 17 May 2017; accepted 20 July 2017; published online xx xx xx)

Conversational speech produced in noise can be characterised by increases in intelligibility relative to such speech produced in quiet. Listening difficulty (LD) is a metric that can be used to evaluate speech transmission performance more sensitively than intelligibility scores in situations in which performance is likely to be high. The objectives of the present study were to evaluate the LD of speech produced in different noise and style conditions, to evaluate the spectral and durational speech modifications associated with these conditions, and to determine whether any of the spectral and durational parameters predicted LD. Nineteen subjects were instructed to speak at normal and loud volumes in the presence of background noise at 40.5 dB(A) and babble noise at 61 dB(A). The speech signals were amplitude-normalised, combined with pink noise to obtain a signal-to-noise ratio of  $-6$  dB, and presented to twenty raters who judged their LD. Vowel duration, fundamental frequency and the proportion of the spectral energy in high vs low frequencies increased with the noise level within both styles. LD was lowest when the speech was produced in the presence of high level noise and at a loud volume, indicating improved intelligibility. Spectrum balance was observed to predict LD. © 2017 Acoustical Society of America.

[<http://dx.doi.org/10.1121/1.4997906>]

[MAH]

Pages: 1–10

## I. INTRODUCTION

Talkers modify their speech in the presence of noise to maintain a level that is sufficient for communication. The Lombard effect (Lombard, 1911) is the involuntary tendency to increase the level of speech in the presence of noise. In noisy environments, speakers commonly increase not only their vocal intensity but also their fundamental frequency ( $f_o$ ; e.g., Junqua, 1993; Van Summers *et al.*, 1988), their first vowel formant ( $F1$ ; e.g., Boril and Pollak, 2005; Kadiri, 1998), and the energy in the spectrum between 1 and 4 kHz relative to the energy below 1 kHz, resulting in an increase in the *spectrum balance* (hereafter SB; e.g., Stanton *et al.*, 1988; Ternström *et al.*, 2006; Krause and Braida, 2004, 2009; Lu and Cooke, 2009). As the speech level and, therefore, the vocal effort level increases, the spectrum flattens (e.g., Nordenberg and Sundberg, 2004; Ternström *et al.*, 2006). When this level increase co-occurs with an increase in glottal flow and, hence, subglottal pressure,  $f_o$  typically rises, while  $F1$  rises with more jaw opening (Fant, 1997). Speech produced in noise can also demonstrate changes in segment (especially vowel) duration and/or a slowing of the speech rate (e.g., Fonagy and Fonagy, 1966; Junqua, 1993; Krause and Braida, 2004).

So-called Lombard speech is more intelligible than speech produced in quiet environments when presented at equal signal-to-noise (SN) ratios (e.g., Dreher and O'Neill,

1957; Summers *et al.*, 1988; Pittman and Wiley, 2001; Lu and Cooke, 2009). However, it has not yet been resolved which of the speech modifications contribute most strongly or are necessary, either in combination or in isolation, to a gain in intelligibility (whether linguistic or non-linguistic parameters; cf. Cooke and García Lecumberri, 2012). See Cooke *et al.* (2014a) for a review. Relatedly, it has not yet been determined which of these speech modifications predict perceived listening difficulty (LD; e.g., Morimoto *et al.*, 2004). However, an upward shift in the overall spectral “centre of gravity” (CoG) does appear to contribute more to intelligibility than does an  $f_o$  increase (e.g., Hazan and Markham, 2004; Lu and Cooke, 2009; Mayo *et al.*, 2012) or the sorts of durational changes that occur in Lombard speech (Cooke *et al.*, 2014b).

Under some conditions, noise-induced speech modifications may be harmful to intelligibility in quiet conditions. Findings for both shouted speech (e.g., Pickett, 1956; Junqua, 1993) and non-native listeners (Cooke and García Lecumberri, 2012) indicate that the speech level may be increased to preserve audibility to the detriment of phonetic information (Rostolland and Parant, 1975). Junqua (1993) observed variation in the intelligibility of Lombard speech relative to speech produced in quiet, depending on the vocabulary, noise type (white Gaussian or multi-talker) and talker gender. For non-native vs native listeners, Cooke and García Lecumberri (2012) found that Lombard speech may be slightly less intelligible than conversational speech when presented in quiet. However, Lombard speech may provide benefits to both native and non-native listeners by placing important speech information outside of the range

<sup>a)</sup>An earlier version of this study was presented at the 45th Annual Symposium of the Voice Foundation, Philadelphia, PA, June 2016.

<sup>b)</sup>Electronic mail: [s.graetzer@liverpool.ac.uk](mailto:s.graetzer@liverpool.ac.uk)

81 of the energetic masker (see discussion in [García](#)  
82 [Lecumberri et al., 2010](#); [Cooke et al., 2014a](#); [Godoy et al.,](#)  
83 [2014](#); [ISO 226, 2003](#)).

84 Previous research has shown that there are differences  
85 in speech intelligibility between Lombard and “clear” speech  
86 or interlocutor-directed speech, such as speech directed  
87 toward infants, hearing-impaired persons, and non-native  
88 speakers (e.g., [Picheny et al., 1985, 1986](#); [Skowronski and](#)  
89 [Harris, 2006](#); [Wassink et al., 2007](#); [Godoy et al., 2014](#);  
90 [Cooke et al., 2014a](#)). Some clear speech modifications are  
91 enhancements that are dependent on linguistic knowledge,  
92 and therefore favour the native speaker (e.g., [Picheny et al.,](#)  
93 [1986](#); [Bond and Moore, 1994](#); [Bradlow and Bent, 2002](#);  
94 [Hazan and Markham, 2004](#)).

95 The acoustic and durational differences between Lombard  
96 and loud or shouted speech have been considered by Stanton  
97 and colleagues (e.g., [Stanton, 1988](#); [Stanton et al., 1988](#)), and  
98 [Bond and Moore \(1990\)](#), but much remains to be investigated.  
99 [Stanton \(1988\)](#) compared the speech modifications associated  
100 with the change from normal speech to speech produced at  
101 “nominally 10 dB above normal,” to the modifications associ-  
102 ated with the change from normal to Lombard speech (involv-  
103 ing 90 dB of pink noise being emitted into the talker’s ears via  
104 headphones) in the fighter cockpit environment. He noted a  
105 smaller shift in the spectral CoG and *F1* and, typically, a  
106 smaller increase in vowel duration between normal and  
107 Lombard speech than between normal and loud speech,  
108 although there was large inter-speaker variation. [Bond and](#)  
109 [Moore \(1990\)](#) concluded, based on a single speaker’s produc-  
110 tion, that Lombard speech and deliberately loud speech  
111 (involving an instruction concerning imagined speaker-listener  
112 distance) result from the same speech production mechanisms.

113 Intelligibility assessment in speech communication can  
114 be performed by means of objective and subjective methods,  
115 such as by calculating the Speech Transmission Index (STI)  
116 of a transmission channel ([IEC 60268-16, 2011](#)), or by test-  
117 ing with real listeners the percentage of words correctly  
118 understood within a given space (intelligibility scores, or  
119 IS). However, sentence scores of 100% are associated in the  
120 [ISO 9921 \(2003\)](#) standard with a large range of STI values  
121 (0.45–1). A sentence intelligibility score of 100% does not  
122 imply that each word is clearly understood, and there are  
123 many situations in which the speech transmission perfor-  
124 mance cannot be regarded as satisfactory ([ISO 9921, 2003](#);  
125 [Morimoto et al., 2004](#), p. 1609). Additionally, for the same  
126 communication channel, scores can be high while predict-  
127 ability and/or word familiarity are high, but reduce when  
128 words are unpredictable or unfamiliar (e.g., [Kalikow et al.,](#)  
129 [1977](#)). These issues of metric sensitivity in the context of  
130 high performance and highly familiar and/or predictable  
131 speech material can be resolved with the use of a rating scale  
132 concerning how difficult a given listening situation is (e.g.,  
133 [ITU-T P.85, 1994](#); [IEC 60268-16, 2011](#)). LD is a subjective  
134 perception metric developed by [Morimoto, Sato and col-](#)  
135 [leagues](#) for use with highly familiar words that can be used  
136 to evaluate speech transmission performance more accu-  
137 rately and sensitively than IS in situations in which the per-  
138 formance is likely to be high ([Morimoto et al., 2004](#)). It is  
139 designed to minimise the potential confounding effects of

word familiarity and predictability and the extent of higher 140  
cognitive processing. LD ratings using the 0–3 rating system 141  
described by [Sato and colleagues \(Sato et al., 2005\)](#) are 142  
mapped to IS and STI values in [IEC 60268-16E \(2011\)](#). LD 143  
has been used as a complement to IS or the STI in several 144  
publications concerning the transmission of Japanese or 145  
Korean speech (e.g., [Morimoto et al., 2004](#); [Sato et al.,](#) 146  
[2005](#); [Lee and Jeon, 2011](#)). 147

148 While the listening effort scale has been expanded from 148  
5 to up to 13 or more levels in order to avoid floor or ceiling 149  
saturation effects ([ITU-T P.85, 1994](#)), the LD traditionally 150  
has only 4 levels (from not difficult to extremely difficult), 151  
and is defined as the percentage of the total number of 152  
responses that indicates some level of difficulty. The use of 153  
only four levels can lead to an accumulation of values at the 154  
upper bound ([Morimoto et al., 2004](#); [Genta et al., 2013](#)), 155  
while averaging over the total number of responses means 156  
that variability associated with the individual listener’s 157  
responses cannot be modeled. This variability, which can be 158  
high (see, e.g., [Lee and Jeon, 2011](#); [Genta et al., 2013](#)), may 159  
be due to individual differences in cognitive ability or prefer- 160  
ence. Studies of category scale design have indicated that 161  
data quality (e.g., reliability, sensitivity) tends to improve as 162  
the number of answer categories increases (e.g., [Alwin,](#) 163  
[1992](#)). An alternative seven-point scale for rating LD was pro- 164  
posed by [Gover and Bradley \(2007\)](#), and a five-point scale 165  
attempting to address the saturation issue but not the variation 166  
issue was proposed by [Genta et al. \(2013\)](#), who suggested on 167  
the basis of their results that there was a need for alternative 168  
implementations of the method. A ten-point scale LD metric 169  
and statistical approach designed to address both issues of sat- 170  
uration and listener variation is presented in this paper. An 171  
additional contribution of the paper is the use of LD ratings 172  
with first language English speakers and listeners who have 173  
been audiometrically tested for normal hearing. 174

175 The consideration of LD independently from speech intelli- 175  
gibility is particularly important for hearing aid users, young 176  
children, and older listeners. This is because even under condi- 177  
tions of perfect speech intelligibility, adverse conditions such as 178  
background noise can impair memory of spoken items and lis- 179  
tening comprehension (e.g., [Pichora-Fuller, 2003](#)). The literature 180  
indicates that there are many acoustical modifications of speech 181  
associated with ease of listening that may or may not co-occur 182  
with improved intelligibility such as modifications of speech rate 183  
*f<sub>0</sub>*, formant frequencies, and *f<sub>0</sub>* modulation ([Bond and Moore,](#) 184  
[1994](#); [Lu and Cooke, 2009](#); [Cooke et al., 2014a](#)). Decreased LD 185  
is likely to reduce listener fatigue, which may lead to intelligibil- 186  
ity improvements in extended listening tasks ([Lim and](#) 187  
[Oppenheim, 1979](#)). In contexts in which listening is difficult, 188  
acoustic treatments or signal enhancement may be used to 189  
reduce fatigue and improve recall ([Lim and Oppenheim, 1979](#)). 190

191 In summary, there has been much work contributing to 191  
the understanding of speech in noise, and features of clear 192  
speech and interlocutor-directed speech. However, the ques- 193  
tion of which speech modifications inherently improve intel- 194  
ligibility and reduce LD for the normal hearing native 195  
English speaker has not yet been fully resolved. Moreover, 196  
there is a need for further investigation and modification of 197  
the LD metric. 198

199 In the present study, the principal aim was to evaluate  
 200 acoustic and durational modifications that occur in non-  
 201 communicative laboratory speech in noisy environments in two  
 202 different speech styles and to relate them to perceived LD.  
 203 Based on this aim, the primary research question was as fol-  
 204 lows, where speech was produced in babble noise at 61 dB(A)  
 205 or in background noise at 40.5 dB(A): Do any of the spectral or  
 206 durational measures considered— $f_o$  (in semitones),  $f_o$  modula-  
 207 tion, SB, or vowel duration—predict LD ratings when SN-  
 208 ratio =  $-6$  dB? A further aim was to extend previous results  
 209 (Morimoto *et al.*, 2004; Sato *et al.*, 2005) that indicated that  
 210 LD may be a useful measure of rating speech transmission  
 211 when word recognition performance would be high. This new  
 212 work extends the previous work by using first-language normal  
 213 hearing English talkers and listeners and also by including an  
 214 evaluation of which spectral or durational changes predict LD  
 215 ratings for continuous speech produced in noise. In this way,  
 216 the current paper responds to a call for studies examining the  
 217 relationship between the intonational, spectral, and durational  
 218 features of Lombard speech and speech intelligibility (Cooke  
 219 *et al.*, 2014b). Furthermore, while Lu and Cooke (2009) con-  
 220 sidered the relevant contributions of only  $f_o$  and spectrum flat-  
 221 tening parameters, and suggested that durational increases  
 222 might, like spectrum flattening, contribute to intelligibility, in  
 223 the current study,  $f_o$ , spectrum flattening, and a durational  
 224 parameter are considered. It is worth considering whether the  
 225 acoustical parameters that predict speech intelligibility also pre-  
 226 dict listening difficulty. The findings have implications for the  
 227 improvement of communication in noisy environments.

## 228 II. EXPERIMENTAL PROCEDURES

229 This study was conducted with approval from and in  
 230 accordance with the policies of Michigan State University's  
 231 Human Research Protection Program (IRB No. 13-1149).  
 232 Participants were not compensated. MATLAB v2014b and Praat  
 233 v5.4.01 (Boersma and Weenink, 2015) were used for signal  
 234 processing. Post-processing and statistical analysis were con-  
 235 ducted in R v3.1.2 (R Development Core Team, 2016).

### 236 A. Experiment one: Speech assessment

#### 237 1. Subjects and instructions

238 Nineteen native American English speaking subjects  
 239 (nine males, ten females) of between 18 and 29 years of age  
 240 with a mean of 21 years of age and with self-reported normal  
 241 speech and hearing were recruited. The subjects were  
 242 recorded while reading the "Rainbow" passage text in a  
 243 semi-reverberant room (a classroom) in two different styles,  
 244 corresponding to normal and loud voice levels. In the envi-  
 245 ronment of the talker was multi-talker children's babble  
 246 (classroom babble; *high level*) noise and/or (naturally occur-  
 247 ing) background (*low level*) noise, which was primarily  
 248 associated with the heating, ventilation, and air conditioning  
 249 system. The instructions given for the styles were as follows:  
 250 normal: "Speak in your normal voice" and loud: "Imagine  
 251 you are in a classroom and you want to be heard by all of the  
 252 children." Investigators were present in the room, observing  
 253 the talker.

## 2. Room acoustic measurements and pre-processing 254

255 The recording took place in a classroom of dimensions 255  
 256  $5.8\text{ m} \times 6\text{ m} \times 2.7\text{ m}$ . The floor and ceiling were covered by 256  
 257 absorbent material (carpet and absorbent tiles). Room acous- 257  
 258 tic parameters were measured in an unoccupied state without 258  
 259 furniture from the impulse responses (IRs) generated by bal- 259  
 260 loon pops (according to ISO 3382-2, 2008).  $T_{30}$  was derived 260  
 261 by means of the AURORA software suite (Farina, 2010). 261  
 262 The mid-frequency reverberation time was 0.53 s (standard 262  
 263 deviation = 0.04; see Bottalico *et al.*, 2015). 263

264 The background noise was measured in the unoccupied 264  
 265 room using a Head and Torso Simulator (HATS) Kemar 265  
 266 45BB-1 (G.R.A.S., Denmark). The primary noise source 266  
 267 contributing to the level of 40.5 dB(A) in the talker position 267  
 268 was the ventilation system. Given that the level was below 268  
 269 43 dB(A), the level of speech production in the background 269  
 270 noise condition was not affected by the noise (Lazarus, 270  
 271 1986; Bottalico *et al.*, 2017). 271

272 The multi-talker noise was emitted by a directional loud 272  
 273 speaker (Yamaha studio monitor model HS5, Yamaha, Japan) 273  
 274 at a level of 61 dB(A) in the talker position. This level repre- 274  
 275 sents a common noise level (hereafter,  $L_{noise}$ ) generated by 275  
 276 children in a classroom engaged in quiet group work or individ- 276  
 277 ual work with some movement (Shield and Dockrell, 2004). 277  
 278 The spectral maxima in the babble occurred in the 500 Hz and 278  
 279 1 kHz octave bands. Babble noise was emitted by the loud 279  
 280 speaker rather than by headphones to avoid the perturbation of 280  
 281 the talkers' self-monitoring of auditory feedback. Arguably, if 281  
 282 noise is delivered via headphones, the headphones can alter the 282  
 283 talker's perception of their own voice (due to the effects on 283  
 284 both internal and external hearing), and therefore the talker's 284  
 285 voice production (e.g., Garnier and Henrich, 2014). The babble 285  
 286 signal had deep amplitude fluctuations, while the background 286  
 287 noise was stationary. The mean  $f_o$  of the babble was 256 Hz, 287  
 288 which is within the normal range for children (Titze, 2000). In 288  
 289 the babble noise condition, the SN-ratio of the speech signal 289  
 290 (represented by the concatenated voiced segments) and corre- 290  
 291 sponding noise signal as acquired by the head-mounted micro- 291  
 292 phone was estimated at +24 dB on average in the loud style, 292  
 293 and +22 dB on average in the normal style. 293

294 The speech signal was acquired by an omnidirectional 294  
 295 head-mounted microphone (Glottal Enterprises M-80, 295  
 296 Syracuse, NY) placed at a distance of 5 cm from the mouth 296  
 297 (much less than the critical distance; hence, the signal was 297  
 298 associated only with the direct sound of the talker). The 298  
 299 microphone has a fairly flat frequency response  $<4$  kHz, 299  
 300 with a rising frequency response between 4 and 6 kHz, and a 300  
 301 sensitivity of  $-65\text{ dB} \pm 3\text{ dB}$ . The signal was acquired by a 301  
 302 Roland R-05 digital recorder (Hamamatsu, Japan) in 16 bit/ 302  
 303 44.1 kHz WAV format. The microphone line out was con- 303  
 304 nected to a personal computer (PC) via an external sound 304  
 305 board (Scarlett 2i4 Focusrite, High Wycombe, UK). The sig- 305  
 306 nal was recorded with Audacity v2.0.6 with a sampling rate 306  
 307 of 44.1 kHz. Recordings varied in length between 25 and 307  
 308 45 s, depending on the talker. 308

309 Words were manually segmented in Praat. For the vowel 309  
 310 duration analysis, individual vowels were segmented in 310  
 311 Python v3.4 by means of the FAVE-align and HTK toolkits 311

312 and visually inspected for errors. The FAVE-align toolkit is  
 313 an adaptation of the Penn Forced Aligner, which relies on  
 314 hidden Markov modeling (Rosenfelder *et al.*, 2014). Vowels  
 315 were labeled according to the Carnegie Mellon University  
 316 (CMU) Pronouncing Dictionary representations of the rele-  
 317 vant Rainbow passage words.

### 318 3. Vowel duration

319 Normalised vowel duration was calculated by dividing  
 320 each vowel duration in seconds associated with a given sub-  
 321 ject by that subject's mean in the low Lnoise and normal  
 322 style (a presumed baseline value). Due to heteroscedasticity,  
 323 durations were analysed by means of Welch-corrected one-  
 324 way tests for equal means. Speech rate was considered dur-  
 325 ing testing, but was found not to change in a reliable way  
 326 with the level of noise and so is excluded from the analyses.

### 327 4. Fundamental frequency

328  $f_0$  was extracted from the recordings by means of Praat  
 329 at 10 ms intervals. An autocorrelation-based method was  
 330 used with Hanning windows with a length of 0.043 s, a pitch  
 331 floor of 70 Hz, and a pitch ceiling of 400 Hz.  $f_0$  was then con-  
 332 verted to semitones in R with bases for males and females  
 333 equal to their mean  $f_0$  (Hz): 128 Hz for males and 203 Hz for  
 334 females in this case. These base values are representative of  
 335 typical adult males and females, the difference relating pri-  
 336 marily to differences in membranous vocal fold length  
 337 (Titze, 2000, 2011).

### 338 5. Spectrum balance

339 Sound pressure level (SPL) data concerning the same  
 340 talkers and experimental conditions as in the present study  
 341 have been reported in a previous publication (Bottalico  
 342 *et al.*, 2015). In this previous study, concerning a set of  
 343 speech production data of which the present data are a sub-  
 344 set, it was confirmed that SPL increased in speech produced  
 345 in noisy conditions, specifically, unintelligible children's  
 346 babble at 61 dB(A), relative to speech produced in relatively  
 347 quiet conditions [ambient noise at 40.5 dB(A)]. As in  
 348 Bottalico *et al.* (2015), in the present study, MATLAB version  
 349 2014b was used to obtain a time history of overall SPL eval-  
 350 uated at 0.125 s intervals for each reading of the Rainbow  
 351 passage. The average among all the SPL values was com-  
 352 puted per subject and this mean was subtracted from each  
 353 time history value for that subject (termed  $\Delta$ SPL). This  
 354 within-subject centering was performed in order to evaluate  
 355 the variation in the subject's vocal behaviour in the different  
 356 conditions from their typical vocal behaviour. For each sub-  
 357 ject, the relative amplitudes in each octave band were calcu-  
 358 lated in dB, where each data point corresponded to a  
 359 difference between each level measured in dB for a subject  
 360 and the maximum amplitude calculated for that subject  
 361 across noise and style conditions.

362 Spectral analysis was conducted in order to determine  
 363 whether an increase in the SB occurred in high relative to  
 364 low Lnoise. SB, named after the measure of Ternström *et al.*  
 365 (2006; but modified in form), was defined as the energy

difference between the 1–4 kHz and 0–1 kHz regions or  
*bands* (i.e., the mean energy computed for the upper band  
 minus the mean computed for the lower band, in dB). The  
 upper band limits were chosen on the basis of previous stud-  
 ies (e.g., Krause and Braida, 2004, 2009; Garnier and  
 Henrich, 2014). The SB value will usually be negative, as  
 the low frequency region tends to dominate the voice spec-  
 trum. The SB increases when it goes from more to less nega-  
 tive and, thus, becomes less steep (or in other words, the  
 spectrum becomes more flat). The claim is that in intelligible  
 speech produced by normal talkers, the energy difference  
 between the lower and the upper bands becomes smaller,  
 resulting in an increase in the SB. However, as discussed  
 previously, this difference can also be affected by the speech  
 level and  $f_0$ . SB as here defined relates to the  $\alpha$  ratio measure  
 (but with the negative rather than the positive sign and an  
 upper limit of 4 kHz rather than 5 kHz; see, e.g., Sundberg  
 and Nordenberg, 2006).

In order to measure possible measurement bias due to  
 any babble noise in the signal acquired by the head-mounted  
 microphone, the difference in SB with and without the arti-  
 ficial babble noise for the same speech material was evaluated  
 with a HATS. The same speech material recorded in the  
 same experimental conditions was emitted from the mouth  
 simulator, with and without babble noise being emitted by  
 the loud speaker. The average difference in the SB with and  
 without the babble noise was equal to  $0.12 \pm 1.14$  dB. A  
 paired sample *t*-test indicated that this difference was negli-  
 gible [ $t = -0.67$ , degrees of freedom (df) = 49,  $p = 0.51$ ].

The concatenated words (i.e., the sentences with silen-  
 ces between words removed) produced by each talker in  
 each condition were subjected to long term average spectrum  
 (LTAS) analysis, also performed in Praat. After fast Fourier  
 analysis, each LTAS was calculated and the SB was derived  
 via the “get slope” function with the lower band limits of 0  
 and 1 kHz, and the upper band limits of 1 and 4 kHz, where  
 the energy is averaged over the concatenated signal in dB,  
 based on the mean power of the signal. When the results  
 were compared with those produced with a lower band of  
 50 Hz–1 kHz, the difference was negligible.

An evaluation of the effects of noise, style, and interac-  
 tions of noise and style, noise and gender, and style and gen-  
 der on the response variable, SB, was conducted by means  
 of a linear mixed effects or LME model (*lme4* and *lmerTest*  
 R packages) fitted by restricted maximum likelihood  
 (REML) with the random effects term of talker. The LME  
 model output includes the estimates of the fixed effects coef-  
 ficients, the standard error (SE) associated with the estimate,  
 the df, the test statistic, *t*, and the *p* value. The Satterthwaite  
 method is used to approximate df and calculate *p* values.

## B. Experiment two: LD assessment

Prior to the LD assessment, 20 native American English  
 speaking listeners (10 males, 10 females), who were aged  
 between 18 and 23 years, with a mean age of 21 years, were  
 audiometrically assessed to ensure normal hearing at  
 $\leq 20$  dB hearing level (HL) between 250 Hz and 6 kHz using

an Orbiter 922 v. 2 audiometer (Madsen Kft., Budapest, Hungary) audiometer in a sound-attenuated booth.

### 1. Room acoustic measurements and pre-processing

There were 152 test stimuli per listener (19 talkers  $\times$  2 speaking styles  $\times$  2 noise conditions  $\times$  2 external auditory feedback conditions, which are not considered here). The stimuli were prepared as follows. A short extract of the Rainbow passage (two sentences in length, which did not include the first or the last phrases in the passage) produced by each talker in each condition was linearly amplitude normalised and combined with pink noise in MATLAB to obtain a SN-ratio of  $-6$  dB. This value is the lowest considered by Sato *et al.* (2005). The onset of noise preceded the onset of the signal by 500 ms. The background Lnoise in the listener position in the booth was 25.1 dB(A), as measured using an NTi Measurements microphone M2211 (class 1 frequency response) and analysed by means of NTi XL2 Audio and Acoustic Analyzer (Schaan, Liechtenstein). LD ratings have been used previously with a specific short speech pattern (Kurusu *et al.*, 2013).

### 2. Testing procedures

The stimuli were presented binaurally via Sennheiser HD205 headphones (Wedemark, Germany) in a pseudo-random order to 20 listeners seated in a sound-attenuated booth. Randomisation on the order of presentation and the recording of LD ratings was obtained via a custom Praat script. The instruction was “rate the level of LD for these sentences on a scale of 1 (not difficult, no effort required) to 10 (very difficult, considerable effort required).” Testing was divided into a training phase (8 stimuli) and a testing phase (152 stimuli), and subjects were able to rest between the 2 halves of the testing phase, to reduce any effects of fatigue. The training phase was included and exposure of all listeners to all conditions was specified, in part, to minimise possible context effects (see Sato *et al.*, 2005). The LD assessment took  $\sim 45$  min. Subjects were required to respond to every stimulus.

In the current study, the discrete subjective LD scale was changed from 1 to 4 (as in the original 2004 version of the metric), in which the percentage of values  $>1$  are taken to represent the LD associated with a given experimental condition (Morimoto *et al.*, 2004) to 1 to 10, for reasons outlined in Sec. II B 1.

### 3. Statistical procedures

A cumulative link mixed model (Laplace approximation; ordinal R package) was run with LD as the response variable and Lnoise, style, their interaction, and interactions of both Lnoise and style with talker gender, with both the listener and the talker as random effects terms. To determine which, if any, of the acoustic and durational parameters predicted LD, a LME model fitted by REML was run with LD as the response variable and SB,  $f_o$  (semitones),  $f_o$  (semitones) standard deviation, and normalised vowel duration as independent variables, with an interaction of  $f_o$  (semitones) and gender, and with talker as the random effects term. In

the case of this model, LD was averaged across listeners per signal. Given that the resolution of 1/10 and the SN-ratio of  $-6$  dB led to the LD metric having good coverage of the measurand range, this response variable could be treated as continuous.

## III. RESULTS

First, the effects of Lnoise and style on spectral and durational speech parameters will be reported. Second, the extent to which any of these parameters predict LD will be discussed.

### A. Experiment one: Speech assessment

#### 1. Vowel duration

Welch-corrected one way tests for equal means indicated that there was an effect of Lnoise [ $F(1,18649) = 134.44$ ,  $p < 0.0001$ ], and gender [ $F(1,18985) = 15.75$ ,  $p < 0.0001$ ] but not style ( $p > 0.1$ ) on normalised vowel duration. This effect of Lnoise held per style and per vowel quality [ $i$ ],  $F(1,1062) = 7.25$ ,  $p < 0.01$ ;  $a$ ],  $F(1,149) = 4.83$ ,  $p < 0.05$ ;  $u$ ],  $F(1,528) = 6.60$ ,  $p < 0.05$ ]. As shown in Fig. 1, vowel durations were longer when the speech was produced in high level than low Lnoise, for both males and females.

#### 2. Fundamental frequency

The mean  $f_o$  increased from 200 to 207 Hz from low to high Lnoise for females, and from 125 to 131 Hz from low to high Lnoise for males. Not only the males' but also the females' mean  $f_o$  remained distant from the mean  $f_o$  of the babble signal (256 Hz).

A LME model was built with  $f_o$  (semitones) as the response variable, and as predictors: Lnoise, style, and interactions of Lnoise and style and noise and gender. Talker was included as a random factor. The low Lnoise, the normal style, and the male gender were chosen as the reference levels. As is shown in Fig. 2,  $f_o$  (in semitones) was higher when speech was

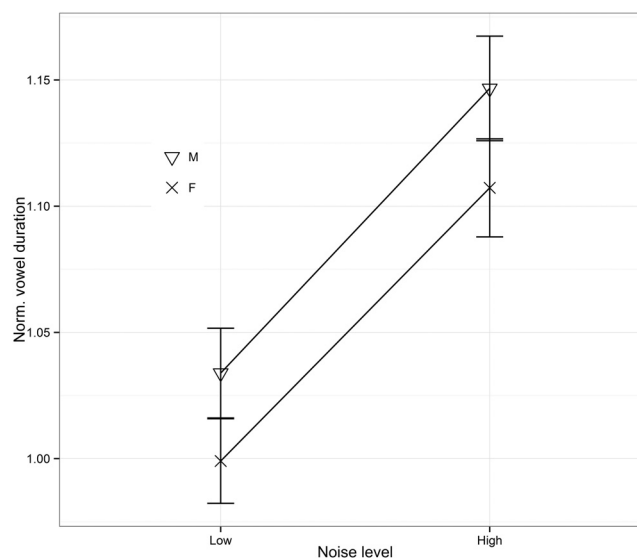


FIG. 1. Normalised vowel durations by Lnoise (x axis) and gender (symbol) condition. Means are shown with 95% confidence intervals.

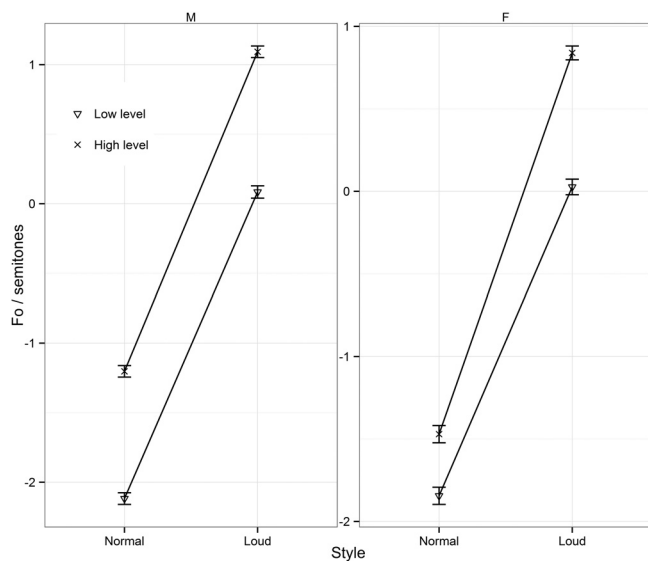


FIG. 2.  $F_0$  in semitones per style (x axis), Noise (symbol) and gender [(left) male, (right) female] condition. Means are shown with 95% confidence intervals.

509 produced in the presence of high Lnoise than low Lnoise  
 510 ( $\hat{\beta} = 0.76$ ,  $SE = 0.03$ ,  $df = 254566$ ,  $t = 29.70$ ,  $p < 0.0001$ ).  $F_0$   
 511 was higher in the loud style than in the normal style ( $\hat{\beta} = 2.09$ ,  
 512  $SE = 0.03$ ,  $df = 254566$ ,  $t = 81.83$ ,  $p < 0.0001$ ). There was an  
 513 interaction between noise and style ( $\hat{\beta} = 0.33$ ,  $SE = 0.03$ ,  
 514  $df = 25466$ ,  $t = 11.52$ ,  $p < 0.0001$ ) such that the effect of noise  
 515 was stronger in the loud style. There was also an interaction  
 516 between style and gender ( $\hat{\beta} = -0.11$ ,  $SE = 0.03$ ,  $df = 25466$ ,  
 517  $t = -3.85$ ,  $p < 0.001$ ), such that males increased their  $f_0$  more  
 518 than females in the loud relative to the normal style.

In the normal style, variation in  $f_0$  (semitones) in the  
 form of standard deviations was slightly increased when  
 speech was produced in the presence of high Lnoise than  
 low Lnoise ( $\hat{\beta} = 0.33$ ,  $SE = 0.14$ ,  $df = 129$ ,  $t = 2.32$ ,  
 $p < 0.05$ ). In the loud style,  $f_0$  variation did not appear to be  
 reliably associated with noise conditions. Variation tended to  
 be lower in the loud style than in the normal style  
 ( $\hat{\beta} = -0.22$ ,  $SE = 0.11$ ,  $df = 129$ ,  $t = -1.96$ ,  $p = 0.05$ ). Very  
 similar results were found when the  $f_0$  values were subjected  
 to outlier detection and removal using the Bonferroni  
 method before analysis, indicating that these results were not  
 due to  $f_0$  artefacts.

### 3. SB

With regard to within-subject normalised overall SPL  
 ( $\Delta$ SPL), in the normal style,  $\Delta$ SPL increased by approximately  
 9 dB from  $-11.74$  dB in low Lnoise to  $-2.70$  dB in high  
 Lnoise. In the loud style,  $\Delta$ SPL increased by approximately  
 4.70 dB from 5.03 dB in low Lnoise to 9.71 dB in high Lnoise.  
 The relative magnitude of spectral energy in the higher frequen-  
 cies was increased in the high Lnoise, as indicated by the relative  
 amplitudes (dB) in each of the seven octave bands (Fig. 3). For  
 males, there tended to be a smaller difference between Lnoise  
 conditions in the loud style than in the normal style, as in the  
 case of the overall  $\Delta$ SPL.

Amplitude variation, measured in terms of the range of  
 the relative amplitude, increased from the low to the high  
 Lnoise in the normal style by 3 dB, and in the loud style for  
 the females by 2 dB, but did not increase in the loud style for  
 males, possibly due to a ceiling effect.

The effects of noise, style, and gender on SB are reported  
 in Table I and shown in Fig. 4. Recall that SB will typically

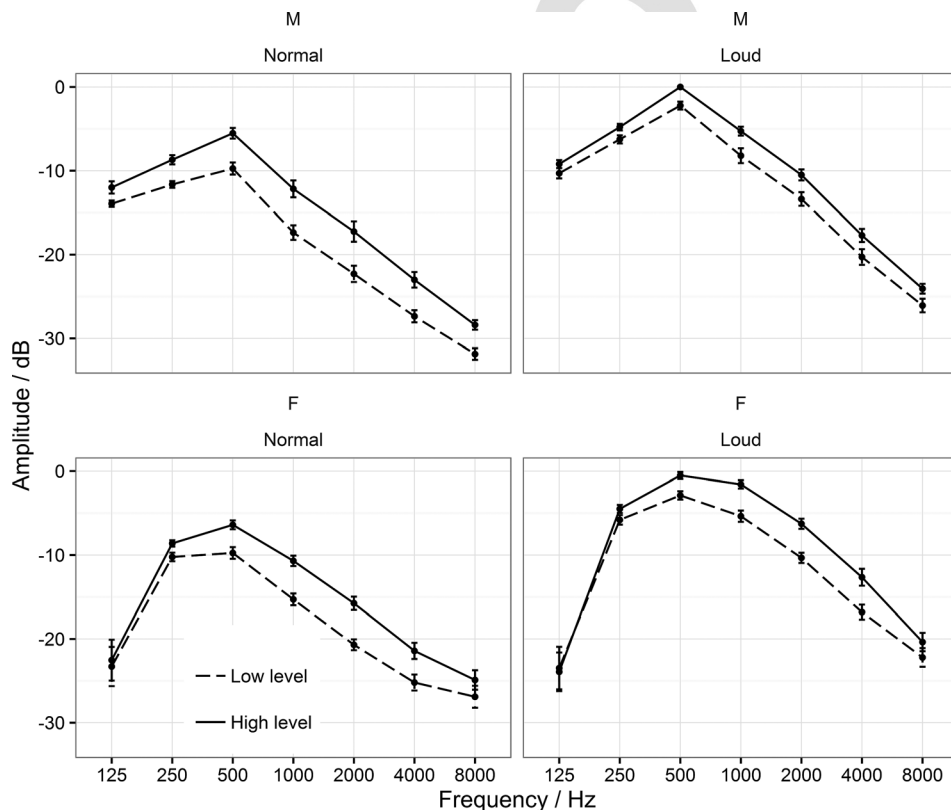


FIG. 3. Relative amplitude (dB) by frequency (Hz), style [(left) loud, (right) normal] and noise level (dashed line, low; solid line, high) for males (upper) and females (lower) with  $\pm 1$  SE.

TABLE I. LME model with the response variable SB and independent variables Lnoise, and style and interaction terms with gender (reference levels: Lnoise, low; style, normal; gender male). Significance codes: \*\*\* < 0.001, \*\* < 0.01, \* < 0.05, “.” < 0.1.

Term	Estimate	SE	df	t
(Intercept)	-17.60	0.55	21	-31.26***
Lnoise high	1.75	0.27	129	6.46***
Style loud	3.31	0.27	129	8.88***
Lnoise high: Style loud	-0.72	0.31	129	-2.34*
Lnoise low: Gender female	0.76	0.75	20	1.02
Lnoise high: Gender female	1.30	0.75	20	1.74
Style loud: Gender female	1.73	0.31	128	5.62***

TABLE II. Cumulative link mixed model (Laplace) output for LD by Lnoise and style and interactions with talker gender (reference levels are Lnoise, low; style, Normal; gender male). Significance codes: \*\*\* < 0.001, \*\* < 0.01, \* < 0.05, “.” < 0.1.

Term	Estimate	SE	z
Lnoise high	-0.86	0.12	-7.40***
Style loud	-1.17	0.12	-10.02***
Lnoise high: Style loud	0.57	0.13	4.37***
Lnoise low: Gender female	-0.61	0.35	-1.74.
Lnoise high: Gender female	-0.61	0.35	-1.73.
Style loud: Gender female	-0.16	0.13	-1.20

550 increase (become less negative) as SPL increases. The LME  
 551 model included Lnoise and style and interactions of Lnoise  
 552 and style, Lnoise and gender, and style and gender, with talker  
 553 as a random effect. Reference levels were low Lnoise, normal  
 554 style, and male gender. When the speech was produced in high  
 555 vs low Lnoise in the normal and in the loud styles, there was  
 556 an increase in SB. In addition, when the speech was produced  
 557 in the loud style vs the normal style in the presence of low  
 558 Lnoise, there was an increase in SB. In the normal style, there  
 559 was a greater difference between Lnoise conditions than in the  
 560 loud style. This result suggests the achievement of an upper  
 561 limit in the high Lnoise, loud style condition. There was also  
 562 an interaction between style and gender: For males there was a  
 563 much smaller difference between the style conditions than for  
 564 the females (Fig. 4). For females,  $f_o$  was moderately positively  
 565 correlated with SB ( $r = 0.52, p < 0.0001$ ). No reliable correla-  
 566 tion was observed for males ( $r = 0.24, p < 0.05$ ).

567 **B. Experiment two: LD assessment**

568 **1. Effects of noise, style, and talker gender on LD**

569 In the speech perception study, 20 listeners evaluated  
 570 their difficulty in listening to the speech produced by the

571 talkers in the 2 noise and 2 style conditions. A cumulative  
 572 link mixed model was fit to LD with the following predic-  
 573 tors: Lnoise and style and interactions of Lnoise and style,  
 574 Lnoise and gender, and style and gender. The model incor-  
 575 porated random effects for talker and listener and for the  
 576 talker listener interaction. The reference levels were low  
 577 Lnoise, normal style, and male gender. As reported in Table  
 578 II and shown in Fig. 5, there was a decrease in LD when the  
 579 speech was produced in high vs low Lnoise in the normal  
 580 style, and when the speech was produced in the loud style vs  
 581 normal style in low Lnoise. There was an interaction of noise  
 582 and style such that the difference in LD between the styles  
 583 was greater in low Lnoise than in high Lnoise. The lowest  
 584 LD scores occurred when speech was produced in high  
 585 Lnoise in the loud style condition.

586 When LD was converted by quartile to a four-point  
 587 scale (as in the original method), the effects of noise and  
 588 style were very similar to those observed in the ten-point  
 589 scale model. In the four-point scale model, the arcsine trans-  
 590 formed proportion of values higher than one (averaged over  
 591 the listeners) was evaluated as the response variable of a  
 592 LME model with noise and style and their interaction as  
 593 independent variables, and talker as a random factor.

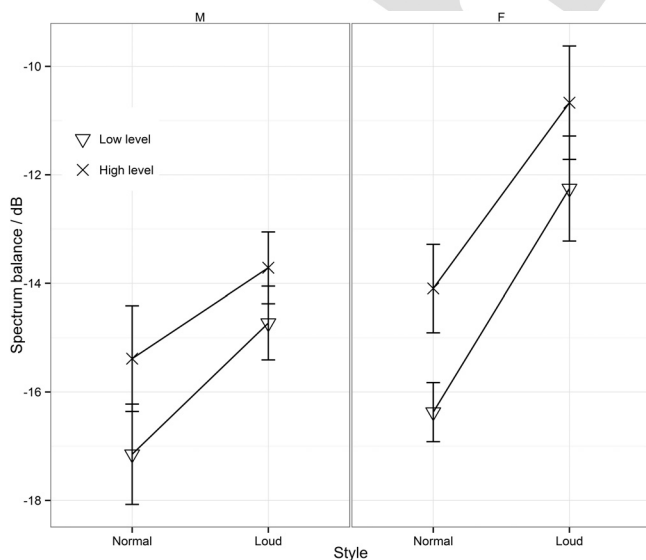


FIG. 4. SB in dB per style (x axis), noise (symbol), and gender [(left) male, (right) female] condition, with means and 95% confidence intervals.

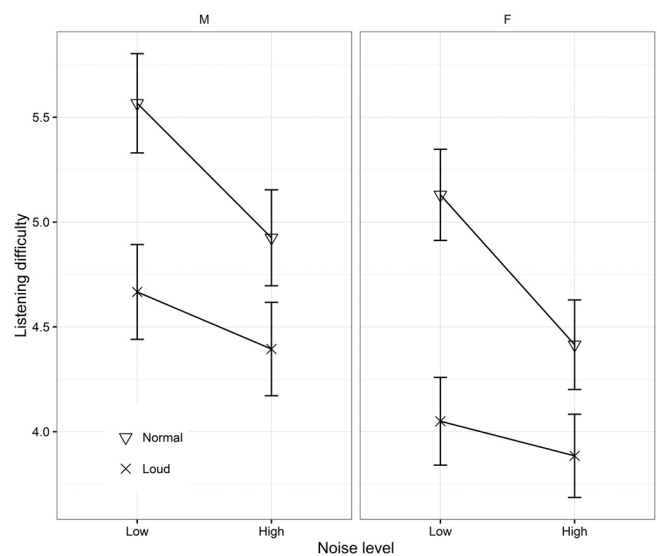


FIG. 5. LD (1, lowest; 10, highest) by Lnoise (x axis) and style (symbol) condition, with means and 95% confidence intervals.

TABLE III. LME model with the response variable LD (averaged over talker) and scaled independent variables: SB,  $f_o$  modulation (semitones), vowel duration, and an interaction of  $f_o$  and talker gender. Significance codes: \*\*\* < 0.001, \*\* < 0.01, \* < 0.05, “.” < 0.1.

Term	Estimate	SE	df	<i>t</i>
(Intercept)	0.05	0.77	135	0.07
SB	-2.9	0.03	82	-10.87***
$F_o$ standard deviation (semitones)	0.11	0.08	50	1.46
Vowel duration	-0.14	0.52	146	0.27
$F_o$ (semitones): Gender male	-0.01	0.04	56	-0.27
$F_o$ (semitones): Gender female	0.13	0.06	71	2.38*

## 594 2. Relationships between speech parameters and LD

595 Models were fitted to determine which of the acoustic  
596 and durational parameters predicted LD. The distribution of  
597 the LD ratings was near normal, with no saturation at the  
598 upper and lower bounds. The results are reported in Table  
599 III. A LME model was run with LD as the response variable  
600 and the acoustical and durational parameters as independent  
601 variables: SB,  $f_o$  modulation (semitones), normalised vowel  
602 duration, and an interaction of  $f_o$  and talker gender. Of the  
603 parameters, only SB reliably predicted LD ( $p < 0.0001$ ).  
604 However, there was a difference in the slope  $f_o$  (semitones)-  
605 LD between females and males such that for females there  
606 was a decrease in LD as  $f_o$  increased ( $p < 0.05$ ). The slope  
607 can be derived from a simple linear regression:

$$y = 4.27 - 0.15f_o + \epsilon. \quad (1)$$

## 608 IV. DISCUSSION

609 This paper reports the use of LD ratings for an identifi-  
610 cation of the speech modifications that predict the transmis-  
611 sion performance of speech produced in noise by first-  
612 language, normal hearing English speakers. In the assess-  
613 ment of the speech parameters in this study, the increase in  
614 vocal intensity in speech produced in noise was found to  
615 co-occur with increases in  $f_o$  and SB, as predicted on the  
616 basis of previous studies (e.g., Van Summers *et al.*, 1988;  
617 Stanton *et al.*, 1988; Junqua, 1993). Arguably, these spec-  
618 tral modifications, which occurred in a non-communicative  
619 context, are primarily associated with the increase in vocal  
620 intensity in the presence of babble noise, but could also  
621 reflect other modifications made to improve audibility for  
622 the talker at his/her own ear (see, e.g., Garnier and Henrich,  
623 2014; Cooke *et al.*, 2014a).

624 In the present study, it was possible to identify effects  
625 of noise within both speech styles. First, for normalised  
626 vowel duration, there was an effect of Lnoise but no  
627 observable effect of style. Additionally, for  $f_o$ , SB, and LD,  
628 there was an additive effect of Lnoise and style. In full, the  
629 effects of noise in the environment of the talker were an  
630 increase in vowel duration, an increase in  $f_o$ , an increase in  
631 the SB, and, in the perception assessment, a decrease in rat-  
632 ings of LD.

633 The results concerning the relationship between dura-  
634 tional changes and LD ratings suggest that while vowel

elongation can increase the amount of acoustic information  
available about vowel quality and neighbouring segment  
identity (see, e.g., Fonagy and Fonagy, 1966), the extent to  
which these changes can improve the intelligibility of speech  
masked by broadband noise appears to depend on other fac-  
tors (Cooke *et al.*, 2014b; Lu and Cooke, 2009). The magni-  
tude of the vowel duration results may reflect the fact that  
the high Lnoise present during speech production was multi-  
talker noise, which is said to degrade the perception of vow-  
els more than consonants (Junqua, 1993). It is interesting  
that there was no reliable effect of style on vowel duration  
for these speakers, despite the observed increase in speech  
level from normal to loud style (cf., e.g., Traunmüller and  
Eriksson, 2000).

Shifts in the spectral energy distribution toward frequen-  
cies between 1 and 4 kHz, i.e., increases in SB, were  
observed to predict LD when the signals were presented to  
listeners at the same SN-ratio. The reported effects of these  
shifts on LD ratings are consistent with the results of Krause  
and Braidá (2004), who found that high frequency spectral  
emphasis contributes to the increased intelligibility of clear  
relative to conversational speech when produced in noise. In  
the present study, it was found that while changes in both  $f_o$   
and spectral energy distribution occur when speech is pro-  
duced in noise, only the latter appears to contribute in a sig-  
nificant way to intelligibility (Lu and Cooke, 2009; Hazan  
and Markham, 2004; Cooke *et al.*, 2014b). The  $f_o$  increase  
may under most conditions merely accompany the increase  
in vocal intensity (Gramming *et al.*, 1988). Lu and Cooke  
(2009) have argued that SB, unlike an upward shift in  $f_o$ , reli-  
ably increases the amount of information available to the lis-  
tener, i.e., the amount of speech information that is out of  
the range of the masker energy. In Cooke's (2006) glimpsing  
model of speech perception in noise, there are more glimpses  
(defined as connected regions in the spectro-temporal rep-  
resentation of the speech time-frequency plane) within which  
speech information is audible. In other words, in the current  
study, an increase in SB provides some release from ener-  
getic masking. In singers, an increase in energy in the  
region of 3 kHz allows the voice to be heard well above an  
orchestra or background noise (Sundberg, 1994) and results  
in an increase in phons and sones (Hunter *et al.*, 2006), due  
to the human ear being particularly sensitive to frequencies  
in this region (see Cooke *et al.*, 2014a; ISO 226, 2003). As  
mentioned previously, Cooke and García Lecumberri  
(2012) have argued that while some linguistic enhance-  
ments may exist in Lombard speech, such as greater vowel  
space dispersion (Cooke and Lu, 2010), these may be out-  
weighed by other changes that in fact reduce intelligibility;  
linguistic enhancements, therefore, appear to have a limited  
role.

With regard to  $f_o$ , in the normal style, a small but reli-  
able increase was observed in  $f_o$  modulation (in semitones)  
in high vs low Lnoise, which has previously been interpreted  
as evidence of an active strategy to improve audibility in  
noise (Garnier and Henrich, 2014). For females, the increase  
in  $f_o$  (semitones) with the increase in Lnoise was larger in  
the loud than in the normal style, despite the effect of noise  
on speech level being larger in the normal style, which may



694 also suggest an active strategy to optimise masker release.  
 695 Further, for the female speakers, a possible explanation of the  
 696 finding for the slope  $f_o$  (semitones)-LD may be that within the  
 697  $f_o$  range of the females ( $\sim 160$ – $270$  Hz), when  $f_o$  increases  
 698 there may be some additional release from energetic masking.  
 699 This is due to the migration of spectral energy into higher  
 700 parts of the spectrum (given the wider spacing of harmonics  
 701 at high  $f_o$  frequencies). This release may be associated with  
 702 the presence of a high level of noise and/or the raised intensity  
 703 of the loud style (see, e.g., Cooke *et al.*, 2014a). Such a rela-  
 704 tionship between  $f_o$  and speech intelligibility for female talk-  
 705 ers may only occur at low SN-ratios (Barker and Cooke,  
 706 2007). Indeed, in the current study, for females but not males,  
 707  $f_o$  was moderately positively correlated with SB.

708 The increase in amplitude variation in high Lnoise rela-  
 709 tive to low Lnoise for some speakers may not be entirely  
 710 related to the increase in vocal intensity, but may also reflect  
 711 the intention of these speakers to improve their intelligibility  
 712 once the SN-ratio can no longer be improved (e.g., Picheny  
 713 *et al.*, 1985, 1986; Ternström *et al.*, 2006).

714 On the basis of the results presented, it can be argued  
 715 that LD ratings are sensitive to changes in the audibility or  
 716 intelligibility of speech in contexts in which the performance  
 717 would be high due, in part, to the predictability of the speech  
 718 material rather than strictly to the SN-ratio or the character-  
 719 istics of the masker. The results are consistent with the find-  
 720 ings of Morimoto *et al.* (2004) that “LD is not always high  
 721 when background noise is present.” (p. 1611) This research  
 722 suggests that an artificial increase in the SB, for example,  
 723 generated by a filter that amplifies frequencies  $> 1$  kHz,  
 724 may reduce LD. Such processing is feasible to implement in  
 725 real-time (Skowronski and Harris, 2006). Thus, signals may  
 726 be enhanced to improve comprehension and recall for young  
 727 children and older listeners. In this paper, a revised LD  
 728 method has been presented that addresses the issues of sat-  
 729 uration at ceiling performance and high listener variability  
 730 that have been reported in the literature.

## 731 V. CONCLUSIONS

732 The objectives of the present study were to evaluate the  
 733 LD of speech produced in different noise and style condi-  
 734 tions, evaluate the spectral and durational speech modifica-  
 735 tions associated with these conditions, and determine  
 736 whether any of the spectral and durational parameters pre-  
 737 dicted LD. It was confirmed that speech produced in high  
 738 level babble noise relative to low level background noise  
 739 was associated with an increased  $f_o$ , increased spectral  
 740 energy between 1 and 4 kHz relative to energy below 1 kHz,  
 741 and increased vowel duration. However, only the proportion  
 742 of high to low spectral energy reliably predicted LD for  
 743 normal-hearing listeners.

744 It should be noted that the speech was acquired in the  
 745 high Lnoise condition in the presence of babble noise; how-  
 746 ever, the effects of noise in the acquired signal on SB itself,  
 747 being the difference between the mean energy of the 1–4 kHz  
 748 band and of the  $< 1$  kHz band, were negligible. In this study,  
 749 the ecological validity of tests in terms of proprioception and  
 750 internal and external auditory feedback was prioritised, as was

unconstrained head movement in the loud style (see, e.g.,  
 Lagier *et al.*, 2010; Garnier and Henrich, 2014).

Further studies are required to evaluate the ten-point scale  
 form of the LD measure. In a future study, the properties of  
 the original form, the revised form, and IS will be compared  
 both for repeated and unique speech material. Moreover, not  
 only the level but also the type of noise in the talker’s environ-  
 ment will be manipulated during communicative tasks, for  
 example, among broadband, speech-shaped, and babble noise,  
 to allow a clear separation of the effects on speech audibility  
 and intelligibility of the level from the type of noise. The type  
 of additive noise used in the listening experiment will also be  
 varied to evaluate how LD ratings and word recognition scores  
 are affected by the properties of the noise masker.

## ACKNOWLEDGMENTS

Thanks are due to the subjects who participated in the  
 study and to members of the audience who were present when  
 an earlier version of this work was presented at the Voice  
 Symposium in June 2016. We also wish to thank the associate  
 editor and two anonymous reviewers for their constructive  
 feedback. This research was primarily supported by the  
 National Institute on Deafness and Other Communication  
 Disorders (NIDCD) of the National Institutes of Health (NIH)  
 under Grant No. R01DC012315. The content is solely the  
 responsibility of the authors and does not necessarily  
 represent the official views of the NIH.

- Alwin, D. F. (1992). “Information transmission in the survey interview:  
 Number of response categories and the reliability of attitude mea-  
 surement,” *Sociol. Methodol.* **22**, 83–118.
- Barker, J., and Cooke, M. (2007). “Modelling speaker intelligibility in  
 noise,” *Speech Commun.* **49**, 402–417.
- Boersma, P., and Weenink, D. (2015). “PRAAT: Doing phonetics by com-  
 puter (version 5.4.01) [computer program],” <http://www.praat.org> (Last  
 viewed 1 May 2015).
- Bond, Z. S., and Moore, T. J. (1990). “A note on Loud and Lombard  
 speech,” in *Proc. 1st International Conference on Spoken Language  
 Processing (ICSLP)*, Kobe, Japan, pp. 969–972.
- Bond, Z. S., and Moore, T. J. (1994). “A note on the acoustic-phonetic char-  
 acteristics of inadvertently clear speech,” *Speech Commun.* **14**, 325–337.
- Boril, H., and Pollak, P. (2005). “Design and collection of Czech Lombard  
 speech database,” in *Proc. Interspeech 2005*, Lisboa, Portugal, pp.  
 1577–1580.
- Bottalico, P., Graetzer, S., and Hunter, E. J. (2015). “Effects of voice style,  
 noise level, and acoustic feedback on objective and subjective eval-  
 uations,” *J. Acoust. Soc. Am.* **138**(6), 498–503.
- Bottalico, P., Ipsaro Passione, I., Graetzer, S., and Hunter, E. J. (2017).  
 “Evaluation of the starting point of the Lombard effect,” *Acta Acust.  
 Acust.* **103**(1), 169–172.
- Bradlow, A., and Bent, T. (2002). “The clear speech effect for non-native  
 listeners,” *J. Acoust. Soc. Am.* **112**, 272–284.
- Cooke, M. (2006). “A glimpsing model of speech perception in noise,”  
*J. Acoust. Soc. Am.* **119**, 1562–1573.
- Cooke, M., and García Lecumberri, M. L. (2012). “The intelligibility of  
 Lombard speech for non-native listeners,” *J. Acoust. Soc. Am.* **132**(2),  
 1120–1129.
- Cooke, M., King, S., Garnier, M., and Aubanel, V. (2014a). “The listening  
 talker: A review of human and algorithmic context-induced modifications  
 of speech,” *Comput. Speech Lang.* **28**, 543–571.
- 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**(4), 2059–2069.
- Cooke, M., Mayo, C., and Villegas, J. (2014b). “The contribution of dura-  
 tional and spectral changes to the Lombard speech intelligibility benefit,”  
*J. Acoust. Soc. Am.* **135**(2), 874–883.

- 816 Dreher, J. J., and O'Neill, J. J. (1957). "Effects of ambient noise on speaker intel- 890  
817 ligibility for words and phrases," *J. Acoust. Soc. Am.* **29**(12), 1320–1323. 891
- 818 Fant, G. (1997). "The voice source in connected speech," *Speech Commun.* 892  
819 **22**, 125–139.
- 820 Farina, A. (2010). "Aurora plug-ins," available at <http://www.aurora-plugins.com> 893  
821 (Last viewed 20 January 2016). 894
- 822 Fonagy, I., and Fonagy, J. (1966). "Sound pressure level and duration," 895  
823 *Phonetica* **15**, 14–21.
- 824 García Lecumberri, M. L., Cooke, M., and Cutler, A. (2010). "Non-native speech 896  
825 perception in adverse conditions: A review," *Speech Commun.* **52**, 864–886.
- 826 Garnier, M., and Henrich, N. (2014). "Speaking in noise: How does the 897  
827 Lombard effect improve acoustic contrasts between speech and ambient 898  
828 noise," *Comput. Speech Lang.* **28**, 580–597.
- 829 Genta, G., Astolfi, A., Bottalico, P., Barbato, G., and Levi, R. (2013). 899  
830 "Management of truncated data in speech transmission evaluation for 900  
831 pupils in classrooms," *Measur. Sci. Rev.* **13**(2), 75–82.
- 832 Godoy, E., Koutsogiannaki, M., and Stylianou, Y. (2014). "Approaching 901  
833 speech intelligibility enhancement with inspiration from Lombard and 902  
834 clear speaking styles," *Comput. Speech Lang.* **28**(2), 629–647.
- 835 Gover, B. N., and Bradley, J. S. (2007). "Comparison of subjective and 903  
836 objective ratings of intelligibility of speech recordings," *Can. Acoust.* 904  
837 **35**(3), 140–141.
- 838 Gramming, P., Sundberg, J., Ternstrom, S., Leanderson, R., and Perkins, W. 905  
839 (1988). "Relationship between changes in voice pitch and loudness," 906  
840 *J. Voice* **2**, 118–126.
- 841 Hazan, V., and Markham, D. (2004). "Acoustic-phonetic correlates of talker 907  
842 intelligibility for adults and children," *J. Acoust. Soc. Am.* **116**(5), 908  
843 3108–3118.
- 844 Hunter, E. J., Švec, J. G., and Titze, I. R. (2006). "Comparison of the pro- 909  
845 duced and perceived voice range profiles in untrained and trained classical 910  
846 singers," *J. Voice* **20**(4), 513–526.
- 847 IEC (2011). 60268-16(E). "Sound system equipment, Part 16: Objective rat- 911  
848 ing of speech intelligibility by speech transmission index" (European 912  
849 Committee for Standardization, Brussels, Belgium, 2011).
- 850 ISO 226:2003(E) (2003). "Acoustics—Normal equal-loudness level contours" 913  
851 (International Organization for Standardization, Geneva, Switzerland).
- 852 ISO 9921:2003E (2003). "Ergonomics: Assessment of speech communication" 914  
853 (International Organization for Standardization, Geneva, Switzerland).
- 854 ISO 3382-2:2008(E) (2008). "Acoustics—Measurement of room acoustic 915  
855 parameters, Part 2: Reverberation time in ordinary rooms" (International 916  
856 Organization for Standardization, Geneva, Switzerland).
- 857 ITU-T P.85 (1994). "A method for subjective performance assessment of the 917  
858 quality of speech voice output devices," ITU-T Recommendation P.85 918  
859 (International Telecommunication Union, Geneva, Switzerland).
- 860 Junqua, J. C. (1993). "The Lombard reflex and its role on human listeners 919  
861 and automatic speech recognizers," *J. Acoust. Soc. Am.* **93**, 510–524.
- 862 Kadiri, N. (1998). "Conséquences d'un environnement bruité sur la produc- 920  
863 tion de la parole" ("Consequences of a noisy environment for speech 921  
864 production"), Ph.D. dissertation, Toulouse University, France.
- 865 Kalikow, D. N., Stevens, K. N., and Elliott, L. L. (1977). "Development of a 922  
866 test of speech intelligibility in noise using sentence materials with con- 923  
867 trolled word predictability," *J. Acoust. Soc. Am.* **61**(5), 1337–1351.
- 868 Krause, J. C., and Braidá, L. D. (2004). "Acoustic properties of naturally 924  
869 produced clear speech at normal speaking rates," *J. Acoust. Soc. Am.* 925  
870 **115**(1), 362–378.
- 871 Krause, J. C., and Braidá, L. D. (2009). "Evaluating the role of spectral and 926  
872 envelope characteristics in the intelligibility advantage of clear speech," 927  
873 *J. Acoust. Soc. Am.* **125**(5), 3346–3357.
- 874 Kurisu, K., Nakamura, S., Yasu, K., and Arai, T. (2013). "Signal to overlap- 928  
875 masking ratio of the broadcasted speech and its listening difficulty: An 929  
876 application for an evaluation tool of sound system tuning," *Acoust. Sci.* 930  
877 *Tech.* **34**(5), 354–355.
- 878 Lagier, A., Vaugoyeau, M., Ghio, A., Legou, T., Giovanni, A., and 931  
879 Assaiante, C. (2010). "Coordination between posture and phonation in 932  
880 vocal effort behaviour," *Folia Phoniatr. Logop.* **62**, 195–202.
- 881 Lazarus, H. (1986). "Prediction of verbal communication in noise—A 933  
882 review: Part 1," *Appl. Acoust.* **19**, 439–463.
- 883 Lee, P. J., and Jeon, J. Y. (2011). "Evaluation of speech transmission in 934  
884 open public spaces affected by combined noises," *J. Acoust. Soc. Am.* 935  
885 **130**(1), 219–227.
- 886 Lim, J. S., and Oppenheim, A. V. (1979). "Enhancement and bandwidth 936  
887 compression of noisy speech," *Proc. IEEE* **67**(12), 1586–1604.
- 888 Lombard, É. (1911). "Le signe de l'élévation de la voix" ("The sign of the ele- 937  
889 vation of the voice"), *Ann. Maladies de L'Oreille Larynx* **37**(2), 101–119.
- Lu, Y., and Cooke, M. (2009). "The contribution of changes in F0 and spectral 940  
891 tilt to increased intelligibility of speech produced in noise," *Speech* 941  
892 *Commun.* **51**, 1253–1262.
- Mayo, C., Aubanel, V., and Cooke, M. (2012). "Effect of prosodic changes 942  
893 on speech intelligibility," in *Proc. 13th Annual Conference of the* 943  
894 *International Speech Communication Association (Interspeech)*. 944
- Morimoto, M., Sato, H., and Kobayashi, M. (2004). "Listening difficulty as 945  
896 a subjective measure for evaluation of speech transmission performance in 946  
897 public spaces," *J. Acoust. Soc. Am.* **116**(3), 1607–1613.
- Nordenberg, M., and Sundberg, J. (2004). "Effect on LTAS of vocal loud- 947  
898 ness variation," in *Speech, Music and Hearing I, Quarterly Progress* 948  
899 *Status Report 45* (Royal Institute of Technology, Stockholm), pp. 93–100. 949
- Picheny, M. A., Durlach, N. I., and Braidá, L. D. (1985). "Speaking clearly 950  
902 for the hard of hearing I: Intelligibility differences between clear and con- 951  
903 versational speech," *J. Speech. Lang. Hear. Res.* **28**, 96–103.
- Picheny, M. A., Durlach, N. I., and Braidá, L. D. (1986). "Speaking clearly 952  
905 for the hard of hearing II: Acoustic characteristics of clear and conversa- 953  
906 tional speech," *J. Speech. Lang. Hear. Res.* **29**, 434–446.
- Pichora-Fuller, M. K. (2003). "Processing speed and timing in aging adults: 954  
908 Psychoacoustics, speech perception, and comprehension," *Int. J. Audiol.* 955  
909 **42**, S59–S67.
- Pickett, J. (1956). "Effects of vocal force on the intelligibility of speech 956  
911 sounds," *J. Acoust. Soc. Am.* **28**(5), 902–905.
- Pittman, A. L., and Wiley, T. L. (2001). "Recognition of speech produced in 957  
913 noise," *J. Speech. Lang. Hear. Res.* **44**, 487–496.
- R Development Core Team (2016). "R: A language and environment for statisti- 958  
915 cal computing," available at <http://www.R-project.org> (Last viewed 1 959  
916 May 2016). 960
- Rosenfelder, I., Fruehwald, J., Evanini, K., Seyfarth, S., Gorman, K., 961  
918 Prichard, H., and Yuan, J. (2014). FAVE (Forced Alignment and Vowel 962  
919 Extraction) Program Suite Version 1.1.3 10.5281/zenodo.9846, 963  
920 available at [http://www.research.ed.ac.uk/portal/en/publications/fave- 964  
921 forced-alignment-and-vowel-extraction-suite-version-113\(bbc2046d-6768- 965  
922 47c5-b574-2987895b0307\).html](http://www.research.ed.ac.uk/portal/en/publications/fave-forced-alignment-and-vowel-extraction-suite-version-113(bbc2046d-6768-47c5-b574-2987895b0307).html).
- Rostolland, D., and Parant, C. (1975). "The intelligibility of a foreign lan- 966  
924 guage in a noisy environment," in *Proc. Symposium Intelligibility of* 967  
925 *Speech*, Nottingham, UK.
- Sato, S., Bradley, J. S., and Morimoto, M. (2005). "Using listening difficulty 968  
927 ratings of conditions for speech communication in rooms," *J. Acoust. Soc.* 969  
928 *Am.* **117**(3), 1157–1167.
- Shield, B., and Dockrell, J. E. (2004). "External and internal noise surveys 970  
930 of London primary schools," *J. Acoust. Soc. Am.* **115**(2), 730–738.
- Skowronski, M. D., and Harris, J. G. (2006). "Applied principles of clear 971  
932 and Lombard speech for automated intelligibility enhancement in noisy 972  
933 environments," *Speech Commun.* **48**, 549–558.
- Stanton, B. J. (1988). "Robust recognition of loud and Lombard speech in the 973  
935 fighter cockpit environment," Ph.D. dissertation, Purdue University, 974  
936 West Lafayette, IN.
- Stanton, B. J., Jamieson, L. H., and Allen, G. D. (1988). "Acoustic-phonetic 975  
938 analysis of loud and Lombard speech in simulated cockpit conditions," in 976  
939 *Proc. 13th International Conference on Acoustics, Speech and Signal* 977  
940 *Processing (ICASSP)*, Vol. 1, pp. 331–334.
- Sundberg, J. (1994). "Perceptual aspects of singing," *J. Voice* **8**(2), 978  
942 106–122.
- Sundberg, J., and Nordenberg, M. (2006). "Effects of vocal loudness varia- 979  
944 tion on spectrum balance as reflected by the alpha measure of long-term- 980  
945 average spectra of speech," *J. Acoust. Soc. Am.* **120**(1), 453–457.
- Ternström, S., Bohman, M., and Södersten, M. (2006). "Loud speech over 981  
947 noise: Some spectral attributes, with gender differences," *J. Acoust. Soc.* 982  
948 *Am.* **119**(3), 1648–1665.
- Titze, I. (2000). *Principles of Voice Production* (National Center for Voice 983  
950 and Speech, Iowa), pp. 1–409.
- Titze, I. (2011). "Vocal fold mass is not a useful quantity for describing F0 984  
952 in vocalization," *J. Speech. Lang. Hear. Res.* **54**(2), 520–522.
- Traunmüller, H., and Eriksson, A. (2000). "Acoustic effects of variation in 985  
954 vocal effort by men, women, and children," *J. Acoust. Soc. Am.* **107**(6), 986  
955 3438–3451.
- Van Summers, W., Pisoni, D. B., Bernacki, R. H., Pedlow, R. I., and Stokes, 987  
957 M. A. (1988). "Effect of noise on speech production: Acoustic and percep- 988  
958 tual analyses," *J. Acoust. Soc. Am.* **84**(3), 917–928.
- Wassink, A. B., Wright, R. A., and Franklin, A. D. (2007). "Intraspeaker 989  
960 variability in vowel production: An investigation of motherese, hyper- 961  
962 speech, and Lombard speech in Jamaican speakers," *J. Phon.* **35**(3), 962  
963 363–379.