



HAL
open science

Effects of linear and nonlinear speech rate changes on speech intelligibility in stationary and fluctuating maskers

Martin Cooke, Vincent Aubanel

► **To cite this version:**

Martin Cooke, Vincent Aubanel. Effects of linear and nonlinear speech rate changes on speech intelligibility in stationary and fluctuating maskers. *Journal of the Acoustical Society of America*, 2017, 141 (6), pp.4126-4135. 10.1121/1.4983826 . hal-01615914v2

HAL Id: hal-01615914

<https://hal.science/hal-01615914v2>

Submitted on 15 Oct 2017

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

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

Effects of linear and nonlinear speech rate changes on speech intelligibility in stationary and fluctuating maskers

Martin Cooke*

Language and Speech Laboratory, Universidad del País Vasco, 01006 Vitoria, Spain

Vincent Aubanel

University of Grenoble Alpes, CNRS, GIPSA-lab, Grenoble, France

(Dated: May 2, 2017)

Abstract

Algorithmic modifications to the durational structure of speech designed to avoid intervals of intense masking lead to increases in intelligibility, but the basis for such gains is not clear. The current study addressed the possibility that the reduced information load produced by speech rate slowing might explain some or all of the benefits of durational modifications. The study also investigated the influence of masker stationarity on the effectiveness of durational changes. Listeners identified keywords in sentences that had undergone linear and nonlinear speech rate changes resulting in overall temporal lengthening in the presence of stationary and fluctuating maskers. Relative to unmodified speech, a slower speech rate produced no intelligibility gains for the stationary masker, suggesting that a reduction in information rate does not underlie intelligibility benefits of durationally-modified speech. However, both linear and nonlinear modifications led to substantial intelligibility increases in fluctuating noise. One possibility is that overall increases in speech duration provide no new phonetic information in stationary masking conditions, but that temporal fluctuations in the background increase the likelihood of glimpsing additional salient speech cues. Alternatively, listeners may have benefitted from an increase in the difference in speech rates between the target and background.

PACS numbers: 43.71.Gv,43.72.Dv

* m.cooke@ikerbasque.org; also at Ikerbasque (Basque Science Foundation).

1 I. INTRODUCTION

2 In speech communication scenarios involving the output of natural or synthetic speech,
3 the likelihood of correct message reception in noisy environments can be improved by modi-
4 fying the speech signal prior to output [e.g., 9, 23, 40, 42, 46, 49, 50, 52]. Such approaches are
5 highly-effective: a recent evaluation of 18 algorithms demonstrated gains equivalent to in-
6 creasing signal-to-noise ratio (SNR) by 5.1 dB for natural speech and by 5.6 dB for synthetic
7 speech [The Hurricane Challenge; 14].

8 Most speech modification techniques operate by reallocating energy in time and frequency
9 under a constant input-output RMS energy constraint. Energy reallocation aims to enhance
10 intelligibility by manipulating the spectro-temporal pattern of local SNR, enabling weaker
11 regions to rise above the masker at the expense of portions of the signal whose local SNR is
12 already sufficiently high. Different approaches have variously transferred energy from voiced
13 to voiceless regions [42], boosted some regions of the spectrum at the expense of others [47],
14 enhanced formants [9], increased the amplitude modulation depth of the mid-frequencies
15 [25], or employed dynamic range compression [7] which has the effect of transferring energy
16 from intense to weaker temporal epochs [41, 51]. [13] provides a review of human and
17 algorithmic speech modifications.

18 An alternative to spectro-temporal energy reallocation is the modification of segment or
19 sub-segmental durations. Altered speaking styles such as Lombard speech [e.g., 24, 36, 45],
20 clear speech [35, 48], speech directed at infants [e.g., 21] and speech produced at a distance
21 [19], exhibit durational changes, usually resulting in slower speech, both overall and at the
22 level of individual speech segments. Many of these forms of speech have been found to be
23 more intelligible than unmodified plain speech [17, 26, 35, 37, 43]. While acoustic changes
24 to features such as intensity, spectral tilt and prosody are present in altered speech styles
25 and may play a role in increased intelligibility, it is natural to consider whether durational
26 changes contribute to the improvement. Durational manipulations can be expected to be
27 particularly effective in the presence of fluctuating maskers where the opportunity arises to
28 shift phonetic information in time to regions where the masking source is less intense.

29 Approaches employing durational modifications are less common than those exploiting
30 spectro-temporal energy reallocation. Of the 14 natural speech modification approaches
31 submitted to the aforementioned Hurricane Challenge [14], only three involved significant

32 durational changes. Tellingly, two of these three produced the largest enhancements in
33 the fluctuating masker condition of the Challenge at moderate and adverse SNRs. Indeed,
34 gains of up to 4.4 dB resulted from an approach [GCReTime; 3] that modified durational
35 information only, indicating that alterations to segment durations alone can be a valuable
36 strategy for maskers containing low-frequency temporal modulations.

37 However, the basis for the intelligibility enhancements produced by durational changes
38 is currently unclear. It is possible that listeners are able to take advantage of the reduced
39 information rate of slower speech rate rather than the intended energetic masking release
40 produced by shifting information in time. Evidence for intelligibility benefits of speech
41 rate slowing is mixed. While studies by Adams and colleagues [1, 2] have demonstrated
42 intelligibility increases for slow speech in masking noise, no such effect was observed un-
43 der conditions of simulated hearing loss by [32] nor when linear and nonlinear durational
44 changes observed in Lombard speech were mapped on to plain speech [15]. Intriguingly,
45 while the latter studies used stationary maskers, the sentence material used by [2] and [1]
46 was mixed with four-talker babble, leading to the possibility that the temporal modulation
47 characteristics of the masker played a role in the different outcomes.

48 One goal of the current study was to determine whether a slower speech rate *per se* con-
49 tributes to the intelligibility increases observed in durationally-modified speech. Keyword
50 scores in utterances that had been linearly-elongated were compared with those for utter-
51 ances whose duration was locally-modified in a way designed to minimise energetic masking.
52 If a reduced information rate is responsible for intelligibility gains of durationally-modified
53 speech, we predict that such gains would be observed in linearly-elongated speech, since
54 durational modifications in this case are independent from masker fluctuations.

55 The current study also addressed the issue of whether the intelligibility of durationally-
56 modified speech is affected by the properties of the masker. Utterances were presented in sta-
57 tionary noise and two forms of fluctuating noise: competing speech, and speech-modulated
58 noise with temporal envelope fluctuations matching those of competing speech.

59 II. EXPERIMENT 1: DURATIONAL MODIFICATIONS IN STATIONARY AND 60 FLUCTUATING MASKERS

61 A. Durational modifications

62 Listeners heard sentences which were either unmodified (PLAIN), linearly-stretched
63 (ELONGATED) or nonlinearly-modified (RETIMED). All durational modifications were carried
64 out using the WSOLA algorithm [16] via a sequence of time-scale factors. In the elongation
65 condition, a constant time-scale factor was used, while in the retiming case the time-scale
66 factor sequence was derived from the GCReTime algorithm described in [3] and summarised
67 below.

68 GCReTime is a general-purpose algorithm that takes a pair of acoustic signals and outputs
69 a retimed version of one of them based on the result of optimising a user-defined criterion
70 operating on a comparison of the two input signals. In the current context, the input
71 to GCReTime is a target speech signal and a masker, and the output is a retimed speech
72 signal which maximises a local distance function whose goal to promote the audibility of
73 information-bearing parts of the speech in the presence of the masker. The distance function
74 is maximised using dynamic programming, the end result being a retiming path which
75 defines a sequence of expansions and contractions of the target speech signal. The process
76 is illustrated in Figure 1. Here, the masker (shown at the top of the figure) is a competing
77 speech signal. The target speech and the modified (retimed) speech are drawn on the left
78 and bottom edges of the figure respectively. The unmodified PLAIN condition corresponds
79 to the diagonal path.

80 The GCReTime local distance function $D(i, j)$ is defined on a grid of points i, j corre-
81 sponding to the i th frame of the target speech signal s and the j th frame of the masker
82 m . The local distance function for all possible pairs (i, j) is a matrix, shown as a grayscale
83 image in Fig. 1, where darker regions depict higher values of the function. The local dis-
84 tance function is composed of two components quantifying (1) the masked audibility of the
85 speech signal in frame i in the presence of the masker at frame j , and (2) the informa-
86 tiveness of the speech signal in the vicinity of frame i . The first of these components is
87 operationalised using the glimpse proportion [12], while the second makes use of cochlear-
88 scaled entropy [CSE; 44]; together these components are reflected in the name ‘GCReTime’.

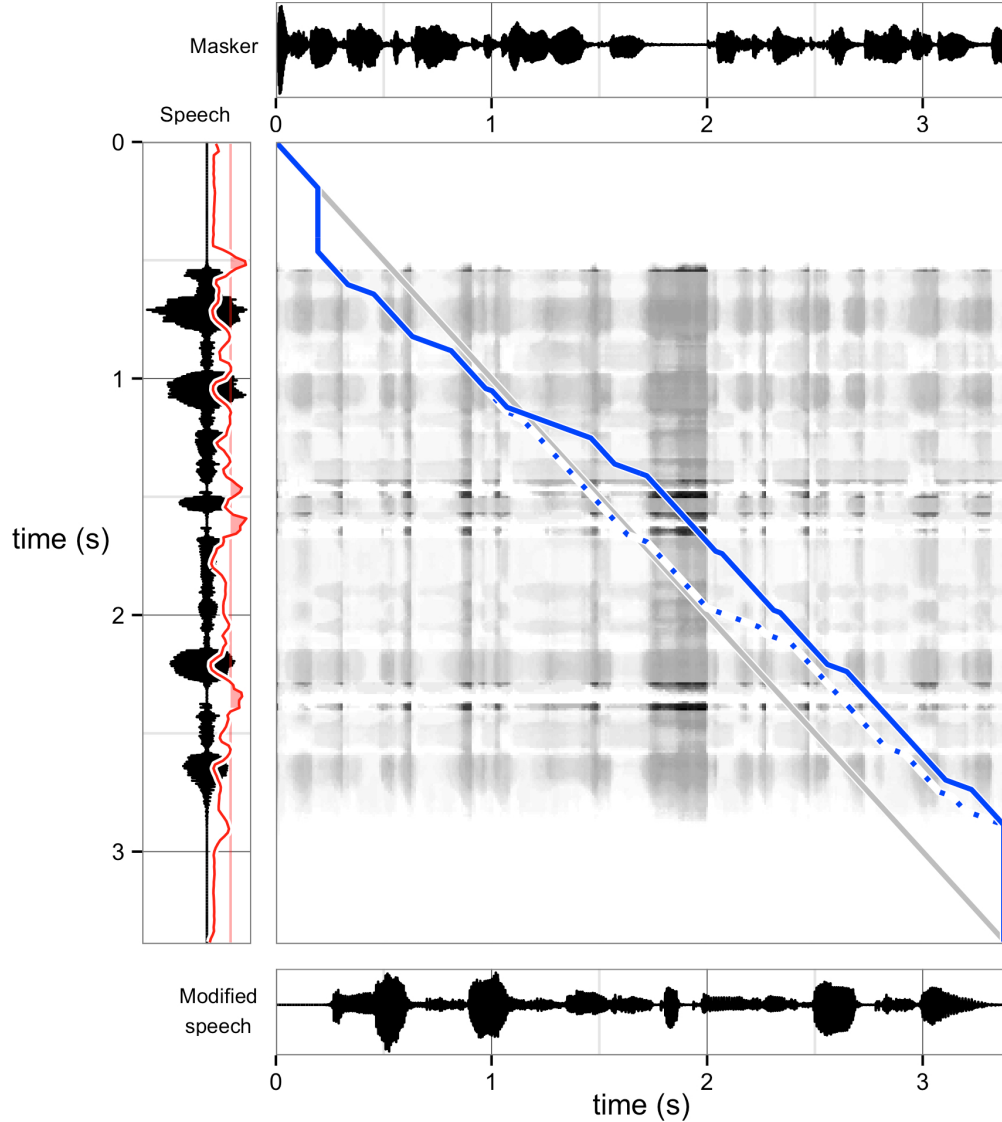


FIG. 1. *An illustration of sentence retiming in the face of a competing speech masker. The grayscale image depicts the value of the local cost function (Eq. A1) for all possible pairs of frames of the target and masker. The solid line shows the minimum cost retiming path using the glimpse proportion (GP) and cochlear-scaled entropy (CSE) components while the dotted line shows the path for the GP component alone. The red curve indicates the value of the CSE weighting defined by Eq. A3; the pink vertical line in the left panel indicates the value of the threshold used to select high-CSE regions (see Appendix). [color online]*

89 CSE captures localised spectral change and has been shown to predict intelligibility better
 90 than consonants, vowels or consonant-vowel/vowel-consonant transitions when tested using
 91 a noise-replacement paradigm [44]. Taken together, these two components ensure that the
 92 distance function takes on high values when the speech signal is not masked and when it
 93 is undergoing a period of rapid change. For example, the dark vertical band in the period

94 immediately preceding the 2 s point in the masker is due to the low level of the masker in
95 that interval, and the darker horizontal strips within this band correspond to those portions
96 of the target speech with a high CSE value. The path that maximises the global distance
97 passes through some of these regions, effectively ensuring that potentially high-information-
98 value transients in the target speech are shifted in time to regions where the masker is less
99 intense.

100 Appendix A describes the computation of the GCReTime local distance function D in
101 more detail.

102 B. Speech and masker materials

103 Utterances were drawn from the phonemically-balanced Sharvard corpus [5] which con-
104 sists of Spanish sentences designed to be equivalent in difficulty to the Harvard sentences
105 [38]. Each Sharvard sentence contains five keywords used for scoring; an example (keywords
106 underlined) is “Llene el frasco de crystal con cola densa” (“Fill the glass flask with thick
107 glue”). Spectrograms of this utterance in each of the three styles PLAIN, ELONGATED and
108 RETIMED are shown in Fig. 2. Renditions of the first 243 Sharvard utterances read by a
109 native Spanish male talker were used in Experiment 1; this figure includes sentences used
110 as practice items.

111 Maskers were constructed using speech material from a native Spanish female talker
112 who read sentences from the Albayzin corpus [30]. Inter-sentence pauses were removed and
113 sentences concatenated to produce a signal of 13.83 minutes duration, sufficient to ensure
114 that no masker fragment was repeated in any speech-plus-masker mixture. Successive non-
115 overlapping fragments from this signal were used for the competing speech masking condition
116 (CS). A speech-shaped noise (SSN) masker was constructed by passing white noise through
117 a filter with a long-term spectrum matching that of the female talker. Each CS fragment
118 had a matched speech-modulated noise (SMN) signal formed by multiplying the short-term
119 temporal envelope of the CS fragment with a portion of the SSN signal selected at random.
120 All speech and masker materials were sampled at 16 kHz.

121 The average PLAIN sentence duration was 2.34 s (s.d. 0.29 s). To allow for overall elon-
122 gation, maskers were constructed to have a duration 0.8 s longer than the target speech
123 utterances they were paired with. Sentences in the ELONGATED and RETIMED conditions

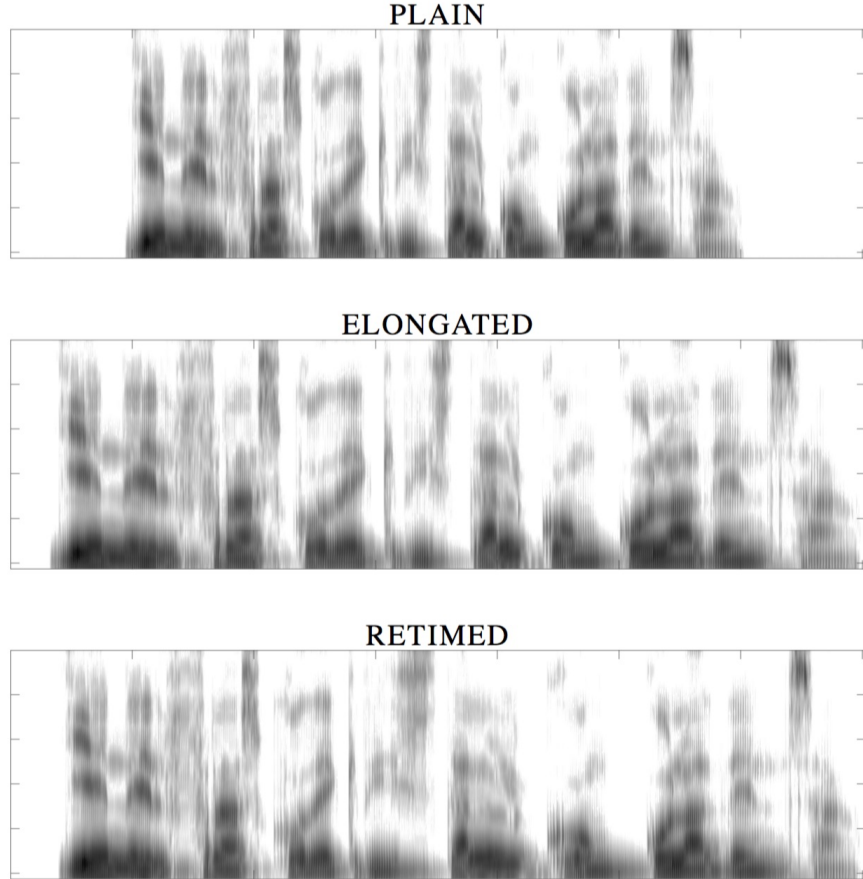


FIG. 2. *Example spectrograms of the utterance “Llene el frasco de cristal con cola densa” in each of the three durational modification conditions.*

124 were 24-55% longer than their PLAIN counterparts (mean 34%, s.d. 4%). Each RETIMED
 125 sentence had a duration that was 97.4-99.2% of the equivalent ELONGATED sentence (mean
 126 98.6%). In Experiment 1, regardless of the masker (CS, SMN or SSN), retiming was carried
 127 out using the CS masker.

128 C. Participants

129 Eighteen native Spanish speakers (10 female) with a mean age of 22.3 years (s.d.=3.8)
 130 took part in the experiment. Speakers were either monolingual in Spanish or bilingual in
 131 Spanish and Basque. All listeners had normal hearing thresholds (< 20 dB HL) in the range
 132 of 125 Hz to 8 kHz, as tested with an Interacoustics AS608 screening audiometer. Listeners
 133 were paid for their participation. Ethics permission was obtained following the University
 134 of the Basque Country ethics procedure.

135 D. Procedure

136 Listeners heard a total of 234 utterances made up of 26 sentences in each of the nine
137 conditions resulting from the combination of the PLAIN, ELONGATED and RETIMED ma-
138 nipulations with the three maskers (CS, SMN and SSN). The SNR for the SSN masking
139 conditions was set to -6.5 dB, a value which led to a 50% keyword score for the male talker
140 in [5]. Since competing speech is a less effective masker than SSN when presented at the
141 same SNR, the SNR for the CS masker was set following pilot tests to -17 dB, while similar
142 tests indicated the need for an intermediate SNR for the SMN case of -12 dB.

143 To avoid sentence subset effects due to possible differing intrinsic intelligibilities of the
144 speech material, the three speech processing conditions (PLAIN, ELONGATED, RETIMED)
145 were applied to the complete set of 234 utterances. Listeners were assigned to subsets of
146 sentences in such a way as to ensure that each sentence in each processing condition was
147 heard the same number of times across listeners, and that each listener heard each sentence
148 exactly once. Speech-plus-noise stimuli were blocked by masker type; within each block
149 listeners heard equal numbers of sentences from each of the three processing conditions in a
150 randomised order. Immediately prior to each block of 78 sentences, listeners responded to
151 3 unscored practice stimuli designed to familiarise them with the type of masker. Practice
152 sentences did not occur elsewhere in the main experiment. Presentation order of the three
153 blocks was balanced across listeners.

154 The listening experiment was conducted in a sound-attenuated studio in the Phonetics
155 Laboratory at the University of the Basque Country, Spain. Speech-plus-noise stimuli were
156 delivered diotically at a presentation level in the range 70.8 - 71.7 dB(A) through Sennheiser
157 HD 380 pro headphones. Listeners received on-screen instructions prior to each block. The
158 experiment ran under computer control using a custom MATLAB program. The experiment
159 was self-paced: following each stimulus presentation participants typed their answer into a
160 text box, after which the next stimulus was presented. On average, listeners required 47
161 minutes (s.d.= 7) to complete the three blocks.

162 E. Postprocessing

163 Listener responses were scored automatically based on the number of keywords identified
164 correctly in each sentence. Vowel stress marks were removed prior to scoring, so that, for
165 example, both ‘mas’ and ‘más’ were considered to be correct responses for the word ‘más’.
166 A total of 130 keywords were scored (26 sentences times 5 keywords per sentence) in each of
167 the nine conditions. Scores were expressed as percentages of keywords identified correctly in
168 each condition. Since none of the scores lay outside the range 23-83%, raw (untransformed)
169 percentages were used in subsequent statistical analyses.

170 F. Results

171 Keyword scores for the PLAIN speech condition were 46.1%, 37.8% and 51.6% for the CS,
172 SMN and SSN maskers respectively.

173 Figure 3 plots changes in scores over the PLAIN baseline for ELONGATED and RETIMED
174 sentences in the three maskers. Elongation led to a small gain in keyword scores of 3.0
175 percentage points (p.p.) in the SSN condition. Substantially larger gains of 8.3 and 9.0
176 p.p. were observed in the two temporally-modulated masking conditions CS and SMN
177 respectively.

178 Retimed speech produced a larger spread of differences over the PLAIN baseline across
179 the three maskers. In stationary noise, retiming was highly detrimental to intelligibility,
180 producing a loss of 14.9 p.p. compared to unmodified speech. For the modulated noise
181 masker (SMN) the gain of 10.3 p.p. was similar to that seen for the ELONGATED condition.
182 However, with a gain of 16.3 p.p., retimed utterances were substantially more intelligible
183 than their elongated counterparts in the competing speech masker.

184 An ANOVA of changes-over-baseline scores with within-subjects factors of modification
185 method (ELONGATED, RETIMED) and masker (CS, SMN, SSN) demonstrated a clear interac-
186 tion in the effect of modifications and maskers [$F(2, 34) = 37.7, p < .001, \eta^2 = 0.28, MSE =$
187 43.1], with significant main effects of both modification type [$F(1, 17) = 10.8, p < .01, \eta^2 =$
188 $0.03, MSE = 20.8$] and masker [$F(2, 34) = 21.2, p < .001, \eta^2 = 0.46, MSE = 164.4$]. Based
189 on a Fisher’s Least Significant Difference (LSD) of 4.4 p.p., gains for ELONGATED speech
190 in the two modulated maskers were equivalent, while ELONGATED speech was statistically-

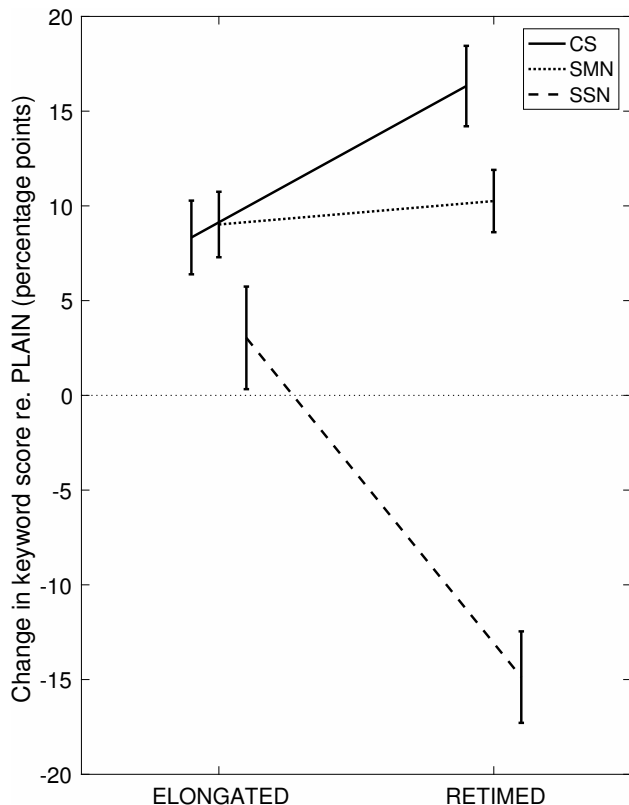


FIG. 3. Changes in mean keywords correct relative to PLAIN speech for ELONGATED and RETIMED utterances in the presence of competing speech (CS), speech-modulated noise (SMN) and speech shaped noise (SSN) maskers. Error bars here and in Fig. 5 represent ± 1 standard error.

191 equivalent to the PLAIN baseline for the SSN masker. Changes in keyword scores for RETIMED
 192 speech were significantly different in the three masking conditions.

193 G. Interim discussion

194 A strategy of retiming speech information by shifting the waveform nonlinearly in time
 195 to attenuate the effect of intense masker epochs has previously been shown to produce
 196 substantial intelligibility gains in the presence of fluctuating maskers [3]. Experiment 1
 197 confirms the effectiveness of algorithmic speech retiming, and extends this finding to speech
 198 material in a different language: the 16.3 p.p. gain produced for Spanish sentences in the
 199 RETIMED condition of the current study in the CS masker condition at an SNR of -17 dB is
 200 consistent with the improvements of 16 and 18 p.p. observed in [3] for English sentences at
 201 SNRs of -14 and -21 dB in the equivalent masking condition of that study.

202 Elongation of speech had a negligible impact on intelligibility for stationary maskers,

203 suggesting that a slower speech rate in itself is not responsible for the gains observed when
204 speech is retimed. In contrast, for fluctuating maskers, elongation led to a clear increase
205 in intelligibility. This finding goes some way to explaining the discrepancies among earlier
206 studies on the effectiveness of a slower speech rate in noise. As noted in the Introduction,
207 while [32] and [15] failed to find an intelligibility benefit of slower speech when presented
208 in a stationary masker, [1] reported a beneficial effect of a slower speech rate in four-talker
209 babble, a type of masker that shows a greater temporal modulation depth than that of a
210 purely stationary noise. The issue of how a fluctuating masker might promote intelligibility
211 increases for elongated speech is addressed in the General Discussion.

212 One intriguing finding is the observation of substantially larger gains produced by retim-
213 ing in the CS condition than in the SMN condition. This outcome would be unexpected
214 if the gains in a fluctuating masker were derived solely from shifting speech information in
215 time to avoid more intense masker intervals. However, since the glimpse proportion analysis
216 underlying GCReTime operates in the spectro-temporal domain, a retiming path produced
217 by a CS masker is not necessarily the same as that produced in response to a SMN masker
218 in spite of the latter having the temporal modulations of the former. Compared to SMN, CS
219 contains some variation in the spectrum across time due to its spectral fine structure of peaks
220 and dips, and it is possible that the retiming path suggested by the GCReTime algorithm
221 is able to take advantage of the glimpsing opportunities afforded in both the spectral and
222 temporal domains. Consequently, the retiming path for CS may be suboptimal for SMN,
223 and vice versa (see Fig. 4). The fact that a competing speech signal was used for retiming
224 in Experiment 1 may have favoured retimed speech when presented in a competing speech
225 masker. To test this hypothesis, a second experiment examined the role of the retiming
226 masker using matched and mismatched retiming maskers.

227 **III. EXPERIMENT 2: ROLE OF THE RETIMING MASKER**

228 **A. Listeners**

229 A new cohort of twenty-one normal hearing paid native Spanish speakers (16 female) with
230 a mean age of 20.0 years (s.d.=1.5) and the same profile as the participants of Experiment 1
231 took part in Experiment 2. Results from one participant who treated the competing speech

232 masker as the target in a number of conditions were excluded.

233 **B. Materials and methods**

234 Raw speech materials, maskers and SNRs were the same as those used in Experiment 1.
235 The ELONGATED condition was not tested. Instead, listeners heard utterances in unmodified
236 PLAIN form and in two distinct retiming conditions. In one retiming condition (CS RETIMED)
237 the durational modifications were based on GCReTime using a CS masker, while in the other
238 retiming condition (SMN RETIMED) the modifications results from the counterpart SMN
239 masker. In this way, listeners heard sentences retimed by a matched or unmatched masker.
240 The nine experimental conditions (3 modifications x 3 maskers) were presented to listeners
241 using the blocking and balancing procedure of Experiment 1 as described in Sec. II D.

242 Figure 4 illustrates the matched/mismatched retiming procedure for the case where re-
243 timed speech was presented in the SMN masker. A comparison of the speech retimed by the
244 SMN masker (matched condition, second panel) and that retimed by the CS masker (mis-
245 matched condition, fourth panel) shows that although the SMN masker is derived from the
246 CS masker, the retimed speech in the matched and mismatched conditions display different
247 temporal structures.

248 **C. Results**

249 Mean keywords correct scores for the PLAIN speech condition were 49.6%, 42.5% and
250 55.8% for the CS, SMN and SSN presentation maskers respectively. Figure 5 plots changes
251 in scores over the PLAIN baseline for sentences retimed using the CS and SMN maskers
252 for each of the three presentation maskers. Changes over baseline for the CS RETIMED
253 conditions were similar to those observed in the equivalent conditions of Experiment 1 (CS,
254 16.3 vs. 16.5 p.p.; SMN 10.2 vs. 10.8; SSN -14.9 vs -12.5), confirming the findings of the
255 first experiment with a different listener cohort.

256 A within-subjects ANOVA with factors of retiming masker and presentation masker for
257 the two fluctuating masking conditions (CS and SMN) indicated no main effect of either
258 factor, but revealed a significant interaction between the two factors [$F(1, 19) = 6.96, p <$
259 $.05, \eta^2 = 0.048, MSE = 32.9$]. Post-hoc η^2 tests based on a Fisher's LSD of 3.80 percentage

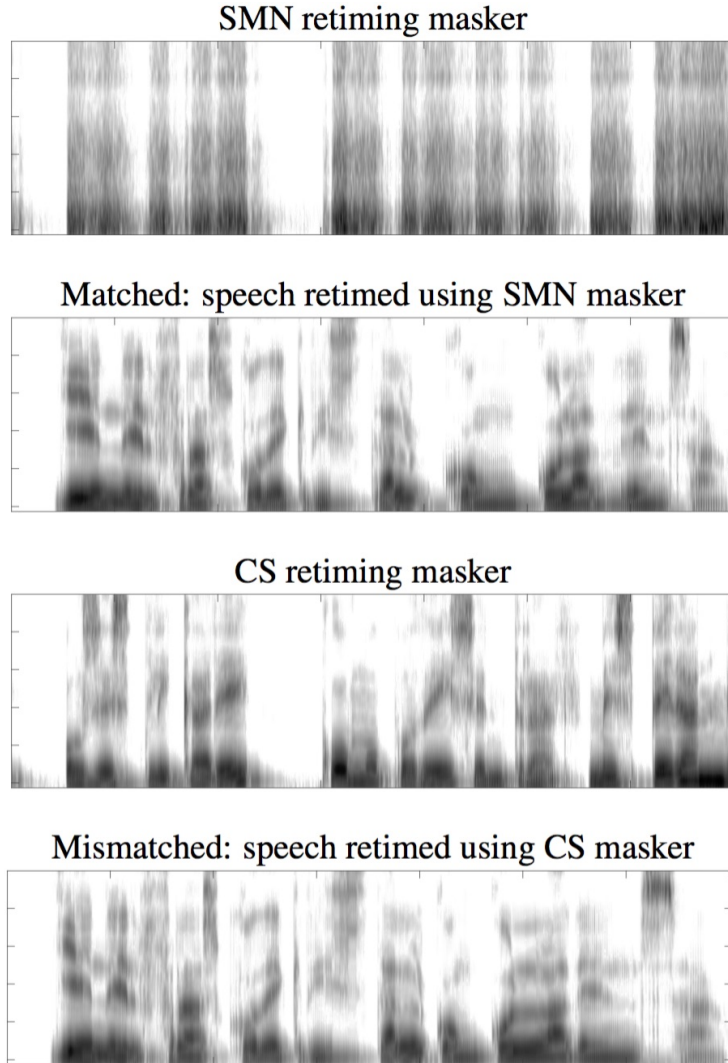


FIG. 4. *Matched and mismatched retiming for the utterance shown in Fig. 2. The top panel shows the SMN masker used to produce the retiming path which results in the utterance shown in the second panel (SMN RETIMED). The third panel shows the CS masker used for retiming which results in the utterance shown in the lower panel (CS RETIMED).*

260 points indicate that gains for CS-based retiming was more effective in a matched CS masker
 261 (16.5 p.p.) than in a mismatched SMN masker (10.9 p.p.). However, there was no benefit
 262 of matched masker type for SMN-based retiming, with similar gains of 13.3 and 12.2 p.p. in
 263 the matched and unmatched conditions respectively. Critically, CS RETIMED speech led to
 264 higher gains than SMN RETIMED speech when presented in a CS masker, suggesting that the
 265 specific details of the retiming path are important. Gains in the two matched conditions (i.e.,
 266 CS RETIMED in CS masker and SMN RETIMED in SMN masker) did not differ statistically,
 267 the difference of 3.2 p.p. falling short of the critical LSD value.

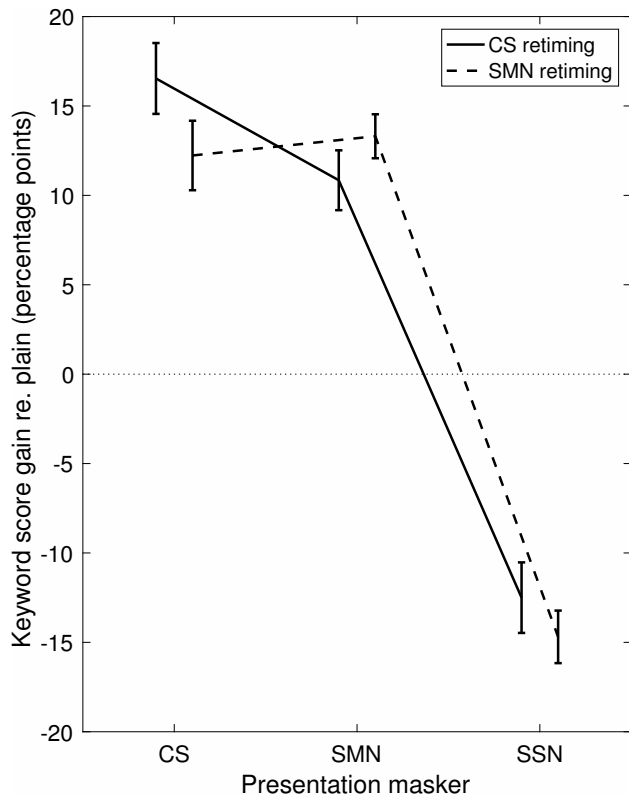


FIG. 5. Changes in keywords correct scores relative to unmodified plain speech for utterances retimed using the CS or SMN masker, for the three presentation maskers.

268 **D. Interim discussion**

269 Experiment 2 demonstrated that the benefits of retiming are affected by the relationship
 270 between the retiming masker and the presentation masker in the case of the CS masker
 271 but not for the SMN masker. This outcome suggests that there is a limit to the benefits of
 272 retiming for a temporally-modulated noise masker, while for CS there may be both temporal
 273 and spectral opportunities which are taken into account by the energetic masking model
 274 underlying the glimpse proportion calculation.

275 **IV. GENERAL DISCUSSION**

276 Experiment 1 addressed the primary research question of the current study by measuring
 277 the extent to which intelligibility gains are present for speech that is linearly elongated to
 278 generate the same average speech rate as that produced by retiming. The absence of a
 279 benefit of elongated speech in the presence of a stationary speech-shaped masker appears to

280 rule out reduced speech rate *per se* as a contributory factor.

281 Nevertheless, elongation led to significant increases in keyword scores in fluctuating
282 maskers, demonstrating that the intelligibility benefits of retiming are not entirely due to
283 the deliberate noise avoidance encapsulated in the GCReTime algorithm. This outcome sug-
284 gests that reducing speech rate can be a very effective strategy for increasing intelligibility
285 in real-life situations characterised by non-stationary sources of noise.

286 There are several ways in which a temporally-fluctuating masker might promote intel-
287 ligibility gains for slowed speech while a stationary masker does not. One possibility is
288 that the regions of elongated speech which escape masking by a stationary noise provide
289 no new phonetic information. The upper two panels of Fig. 6 depict glimpses of speech in
290 the presence of the SSN masker for the example utterance, both unmodified and elongated.
291 It is clear that while small differences in putative glimpses exist due to fluctuations in the
292 speech-shaped noise, the nature of the available information is largely identical in the PLAIN
293 and ELONGATED conditions: the glimpses are simply elongated. For example, the phoneme
294 /k/ in ‘frasco’ is devoid of glimpses in both cases, and the /s/ in the same word conveys the
295 same information in the two cases. In contrast, for the competing speech masker (lower pan-
296 els), temporal fluctuations in the masker increase the likelihood of observing new phonetic
297 information in elongated speech. For example, in the PLAIN speech, there is a paucity of
298 critical low frequency information to indicate the identity of the vowel /e/ in ‘densa’, while
299 such information is present in the ELONGATED version. Of course, while some information is
300 gained in this way, other regions of the signal are likely to be masked with a commensurate
301 loss of information. However, we speculate that since the overall signal duration is increased
302 in the ELONGATED case, so is the net amount of phonetic information.

303 An alternative explanation for the observed gains in fluctuating maskers arises from the
304 possibility that listeners are better able to separate target and background speech due to
305 speech rate differences between the target and masker, overcoming a potential source of
306 informational masking. This notion is supported by a study by [20] in which a cohort of
307 young normal hearing listeners recognised more words in time-compressed sentences when
308 the compression ratio did not match that of 12-talker background babble. The ELONGATED
309 condition of the current study does indeed lead to an increase in speech rate differences
310 between the target and background: the PLAIN speech was articulated at a rate of 5.8 vow-
311 els/s, comparable to the 6.6 vowels/s of the competing speech masker. Speech rate slowing

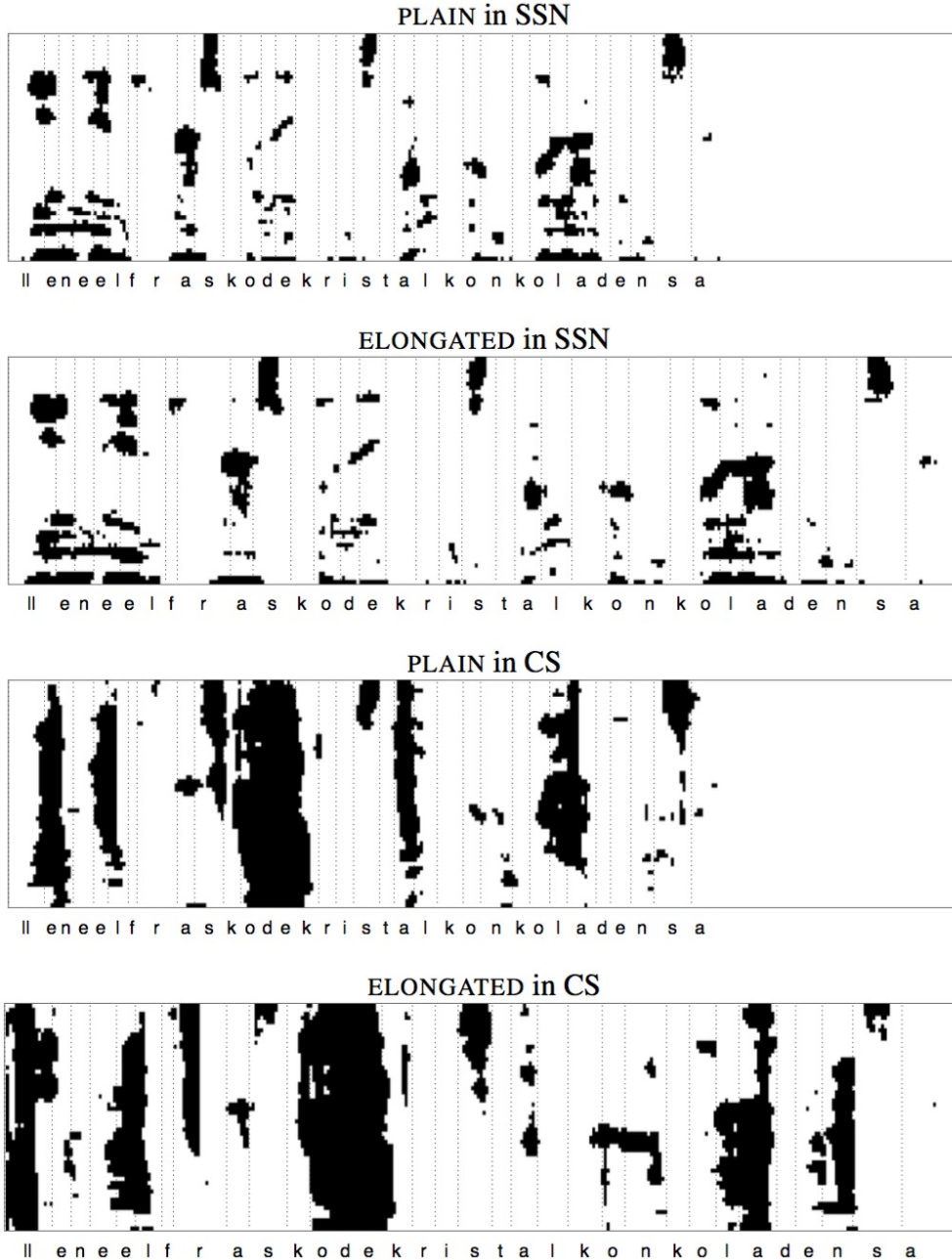


FIG. 6. *Regions of a target speech utterance of the phrase “Llene el frasco de cristal con cola densa” which are deemed to escape energetic masking according to a glimpsing model [12], for PLAIN and ELONGATED speech in the presence of a stationary masker (top two panels) and a fluctuating masker (lower two panels). A broad phoneme-level transcription is provided in each case.*

312 in the ELONGATED condition reduced the average speech rate to 4.4 vowels/s, increasing the
 313 target-background speech rate difference. Earlier studies with listeners [28] or models [8]
 314 have demonstrated sensitivity to differences in speech interruption rate and speech modula-
 315 tion rate. Further studies controlling for speech rate differences are needed to rule out their

316 possible contribution to the observed gains in the current study. However, since gains were
317 also observed for the non-speech masker, it appears necessary to invoke a more generalised
318 notion of temporal modulation rate differences that go beyond speech-on-speech informa-
319 tional masking, which in this case is likely to have a relatively small effect since the gender
320 of the target and masking talkers differed [10]. A related possibility is that the presence
321 of modulation in the masker itself imposes a cognitive burden on listeners: a slower rate
322 of information transmission via a lower speech rate may be beneficial in reducing listening
323 effort.

324 A striking and consistent outcome observed in both experiments was the substantial
325 loss of intelligibility (amounting to 13-15 percentage points) that occurred when speech
326 was nonlinearly-retimed in stationary noise, contrasting with no loss for linearly-elongated
327 speech in the presence of the stationary masker. In the GCReTime algorithm, retiming takes
328 no account of anything other than the temporal relationship between speech and masker
329 (the glimpsing component), and speech dynamics (the CSE component). Consequently,
330 properties such as segment duration and the local speech rate of unmodified speech are
331 not preserved by the algorithm. Such distortions are likely both to confound listeners'
332 expectations of when salient information is going to occur and to diminish the effectiveness of
333 contextual cues that depend on relative durations. For Spanish, changes in relative segment
334 durations induced by GCReTime may have interfered with phonological cues [e.g., 27] or
335 syllabification [e.g., 22]. Intriguingly, such a tradeoff between modifications that overcome
336 masking and those which preserve phonological integrity is also seen in naturally-produced
337 speech in noise. [39] measured durational (vowel lengthening) cues to the voicing distinction
338 in English plosives in plain and Lombard speech, finding a reduced contrast in the latter case.
339 The benefits of retiming in fluctuating maskers presumably reflect a net effect of masking
340 release and durational distortion, suggesting that even larger benefits in noise are realisable
341 if the phonological impact of durational modifications can be minimised.

342 Distortions to the target speech might also have contributed to the observed differences
343 in effectiveness of retiming in the presence of competing speech maskers and temporally-
344 modulated noise, even when the influence of a matching or mismatching retiming masker
345 was controlled for, as in Experiment 2. In a non-informational masker such as modulated
346 noise, listeners' attentional focus is presumably directed to the target speech alone, and any
347 departures from expected phonological forms may be noticeable, and potentially lead to

348 the consideration of additional competitor words. In contrast, when the masker itself con-
349 tains speech, it is conceivable that the cognitive burden imposed by foreground-background
350 separation precludes a more detailed analysis of the target signal, or a mis-attribution of
351 retiming-based distortions to the competing speech signal. Another possibility is that any
352 gains due to retiming outweigh losses due to mistiming of phonological features.

353 The outcome of the current study points to the potential of durational changes as a
354 mechanism for improving intelligibility in noise, but also highlights the need to take the
355 temporal properties of the masker into account, given the deficits resulting from the retiming
356 method in the presence of stationary noise. The finding that gains are possible merely
357 by elongating the speech signal in fluctuating maskers suggests that speech rate slowing
358 could be a component of a simple practical strategy for boosting intelligibility. As noted in
359 the Introduction, modified duration is not by any means the sole manifestation of natural
360 ‘altered’ speaking styles, and spectral factors in particular are known to have a sizeable
361 influence on intelligibility [15]. Spectral and durational changes are orthogonal to a large
362 extent e.g., changes to properties such as spectral tilt can be imposed independently of
363 durational changes.

364 As is typically the case when targetting the 50% correct response rate with a normal-
365 hearing adult population, all testing was done at negative SNRs. Further work is needed to
366 measure the efficacy of a slower speech rate at more realistic SNRs [31], as such environments
367 have also been found to induce slower rate of speech in talkers [4]. The benefits observed in
368 the current study of nonlinear retiming at negative SNRs may be reduced at higher SNRs;
369 lower than expected benefits for a fluctuating masker advantage in comparison to stationary
370 noise have consistently been observed at positive SNRs [6, 18, 33].

371 V. CONCLUSIONS

372 (i) Reductions in speech rate resulting from linear elongation of the speech signal did
373 not lead to intelligibility increases (nor did they disrupt intelligibility) for sentences in the
374 presence of stationary speech-shaped maskers, suggesting that intelligibility gains seen in
375 durationally-modified speech were not due to the reduction in information rate that accom-
376 panies slower speech.

377 (ii) However, identical elongations produced significant intelligibility increases in fluctu-

378 ating maskers. One explanation is that while elongation in stationary maskers produces
 379 no new speech information, the altered pattern of glimpses in fluctuating maskers leads to
 380 the unmasking of new phonetic cues. An alternative is that slower speech enables listeners
 381 to separate target speech from the background due to greater differences in speech rate,
 382 or reduces the cognitive burden of processing speech in a modulated background. Further
 383 studies are needed to distinguish these possibilities.

384 (iii) Nonlinear durational modifications designed to reduce energetic masking of speech
 385 information led to larger intelligibility gains in competing speech maskers than those pro-
 386 duced by linear elongation in spite of the distortion of phonetic integrity indicated by the
 387 reduced intelligibility of the same modifications in stationary maskers.

388 **ACKNOWLEDGEMENTS**

389 This work was partially funded by the “Listening Talker” project, supported by the Future
 390 and Emerging Technologies (FET) programme within the Seventh Framework Programme
 391 for Research of the European Commission, under FET-Open grant no. 256230. Author VA
 392 also acknowledges support from the FP7 FET Project “Speech Unit(e)s”, grant no. 339152.

393 **Appendix A: Computation of the GCRetime local distance function**

394 The GCRetime local distance function (Eq. A1) is defined for each pair of time frames i
 395 of the speech signal and j of the masker as the product of two terms: (i) glimpse proportion,
 396 $GP(i, j)$, the proportion of the speech signal in frame i glimpsed in the presence of the masker
 397 in frame j (Eq. A2); and (ii) $W_{CSE}(i)$, a weighting term based on the cochlear-scaled entropy
 398 of the speech signal in frame i (Eq. A3):

$$D(i, j) = GP(i, j)W_{CSE}(i) \tag{A1}$$

399 **Glimpse proportion**

400 The glimpse proportion is intended to reflect the local audibility of speech in noise and is
 401 defined as the percentage of spectral regions where the modelled auditory excitation pattern

402 for the target speech exceeds that of the masker:

$$GP(i, j) = \frac{1}{F} \sum_{f=1}^F \mathcal{H}[S_f(i) > M_f(j)] \quad (\text{A2})$$

403 where F is the number of frequency channels, S_f and M_f denote the excitation patterns
 404 of speech and masker in frequency channel f , and $\mathcal{H}(\cdot)$ is the Heaviside unit step function
 405 counting the number of channels where the speech exceeds the masker. Excitation patterns
 406 are derived via a gammatone filterbank [34] using an implementation introduced by Cooke
 407 [11]. The Hilbert envelope of each gammatone filter output is computed and smoothed by a
 408 leaky integrator with a 8 ms time constant [29], downsampled and log-compressed. Here the
 409 gammatone filterbank contained $F = 32$ frequency channels spaced equally on an ERB-rate
 410 scale between 50 Hz and 7500 Hz.

411 Cochlear-scaled entropy

412 In the current study we use the concept of cochlear-scaled entropy [CSE; 44] to identify
 413 spectral regions which are changing most rapidly in order to give them greater weight in
 414 the computation of the local distance function. CSE is implemented as a locally averaged
 415 measure of spectral change across time based on excitation patterns of the target speech
 416 signal:

$$CSE(i) = \sum_{k=-\lambda/2}^{\lambda/2} d(i+k)$$

417 where

$$d^2(t) = \sum_{f=1}^F [S_f(t+1) - S_f(i)]^2$$

418 and λ is the number of frames over which the CSE is computed. Following [44], $\lambda = 5$,
 419 equivalent to 80 ms for the 16 ms time frames used here.

420 The CSE-based weighting is defined as

$$W_{\text{CSE}}(i) = (w - 1) \mathcal{H}[CSE(i) - \beta] + 1 \quad (\text{A3})$$

421 where β is a threshold used to identify high-CSE regions, and w defines the degree of boosting
 422 of the CSE value. Here, values of $\beta = 0.6$ and $w = 3$ were used.

-
- 423 [1] Adams, E. M., Gordon-Hickey, S., Morlas, H., and Moore, R. (2012). Effect of rate-alteration
424 on speech perception in noise in older adults with normal hearing and hearing impairment.
425 *Am. J. Audiol.*, 21(1):22–32.
- 426 [2] Adams, E. M. and Moore, R. E. (2009). Effects of speech rate, background noise, and simulated
427 hearing loss on speech rate judgment and speech intelligibility in young listeners. *J. Am. Acad.*
428 *Audiol.*, 20:28–39.
- 429 [3] Aubanel, V. and Cooke, M. (2013). Information-preserving temporal reallocation of speech in
430 the presence of fluctuating maskers. In *Proc. Interspeech*, pages 3592–3596, Lyon, France.
- 431 [4] Aubanel, V., Cooke, M., Villegas, J., and García Lecumberri, M. L. (2011). Conversing in the
432 presence of a competing conversation: effects on speech production. In *Proc. of Interspeech*,
433 pages 2833–2836, Florence, Italy.
- 434 [5] Aubanel, V., García Lecumberri, M. L., and Cooke, M. (2014). The Sharvard Corpus: A
435 phonemically-balanced Spanish sentence resource for audiology . *Int. J. Audiology*, 53:633–
436 638.
- 437 [6] Bernstein, J. G. and Grant, K. W. (2009). Auditory and auditory-visual intelligibility of
438 speech in fluctuating maskers for normal-hearing and hearing-impaired listeners. *J. Acoust.*
439 *Soc. Am.*, 125:3358–3372.
- 440 [7] Blesser, B. A. (1969). Audio dynamic range compression for minimum perceived distortion.
441 *IEEE Trans. on Audio and Electroacoustics*, 17(1): 22–32.
- 442 [8] Bronkhorst, A. W., Bosman, A. J., and Smoorenburg, G. F. (1993). A model for context
443 effects in speech recognition. *J. Acoust. Soc. Am.*, 93:499–509.
- 444 [9] Brouckxon, H., Verhelst, W., and Schuymer, B. D. (2008). Time and frequency dependent am-
445 plification for speech intelligibility enhancement in noisy environments. In *Proc. Interspeech*,
446 volume 9, pages 557–560.
- 447 [10] Brungart, D. S., Simpson, B. D., Ericson, M. A., and Scott, K. R. (2001). Informational and
448 energetic masking effects in the perception of multiple simultaneous talkers. *J. Acoust. Soc.*
449 *Am.*, 100(5):2527–2538.
- 450 [11] Cooke, M. (1993). *Modelling Auditory Processing and Organisation*. Cambridge University
451 Press.

- 452 [12] Cooke, M. (2006). A glimpsing model of speech perception in noise. *J. Acoust. Soc. Am.*,
453 119(3):1562–1573.
- 454 [13] Cooke, M., King, S., Garnier, M., and Aubanel, V. (2014a). The listening talker: A review
455 of human and algorithmic context-induced modifications of speech. *Computer Speech and*
456 *Language*, 28:543–571.
- 457 [14] Cooke, M., Mayo, C., and Valentini-Botinhao, C. (2013). Intelligibility-enhancing speech
458 modifications: the Hurricane Challenge. In *Proc. Interspeech*, pages 3552–3556.
- 459 [15] Cooke, M., Mayo, C., and Villegas, J. (2014b). The contribution of durational and spectral
460 changes to the Lombard speech intelligibility benefit. *J. Acoust. Soc. Am.*, 135(2):874–883.
- 461 [16] Demol, M., Verhelst, W., Struyve, K., and Verhoeve, P. (2005). Efficient non-uniform time-
462 scaling of speech with WSOLA. In *Int. Conf. on Speech and Computers (SPECOM)*, pages
463 163–166.
- 464 [17] Dreher, J. J. and O’Neill, J. J. (1957). Effects of ambient noise on speaker intelligibility for
465 words and phrases. *J. Acoust. Soc. Am.*, 29(12):1320–1323.
- 466 [18] Freyman, R. L., Balakrishnan, U., and Helfer, K. S. (2008). Spatial release from masking with
467 noise-vocoded speech. *J. Acoust. Soc. Am.*, 124:1627–1637.
- 468 [19] Fux, T., Feng, G., and Zimpfer, V. (2012). Natural-to-shouted voice transformation for
469 distance cues of monosyllabic consonant-vowel-consonant words. *Acta Acust. United Ac.*,
470 98(5):839–843.
- 471 [20] Gordon-Salant, S. and Fitzgibbons, P. J. (2004). Effects of stimulus and noise rate variability
472 on speech perception by younger and older adults. *J. Acoust. Soc. Am.*, 115(4):1808–1817.
- 473 [21] Grieser, D. A. L. and Kuhl, P. K. (1988). Maternal speech to infants in a tonal language:
474 Support for universal prosodic features in motherese. *Dev. Psychol.*, 24(1):14.
- 475 [22] Hualde, J. and Chitoran, I. (2003). Explaining the distribution of hiatus in Spanish and
476 Romanian. In *Proc. Int. Conf. Phonetic Sciences*, pages 1683–1686, Barcelona.
- 477 [23] Jokinen, E., Remes, U., and Alku, P. (2016). The use of read versus conversational Lombard
478 speech in spectral tilt modeling for intelligibility enhancement in near-end noise conditions.
479 In *Proc. Interspeech*, pages 2771–2775.
- 480 [24] Junqua, J.-C. (1993). The Lombard reflex and its role on human listeners and automatic
481 speech recognizers. *J. Acoust. Soc. Am.*, 93(1):510–524.
- 482 [25] Koutsogiannaki, M. and Stylianou, Y. (2016). Modulation enhancement of temporal envelopes

- 483 for increasing speech intelligibility in noise. In *Interspeech 2016*, pages 2508–2512.
- 484 [26] Lu, Y. and Cooke, M. (2008). Speech production modifications produced by competing talkers,
485 babble, and stationary noise. *J. Acoust. Soc. Am.*, 124(5):3261–3275.
- 486 [27] Mendoza, E., Carballo, G., Cruz, A., Fresneda, M. D., Muoz, J., and Marrero, V. (2003).
487 Temporal variability in speech segments of Spanish: context and speaker related differences.
488 *Speech Communication*, 40:431–447.
- 489 [28] Miller, G. A. and Licklider, J. C. (1950). The intelligibility of interrupted speech. *J. Acoust.*
490 *Soc. Am.*, 22:167–173.
- 491 [29] Moore, B. C. J., Glasberg, B. R., Plack, C. J., and Biswas, A. K. (1988). The shape of the
492 ear’s temporal window. *J. Acoust. Soc. Am.*, 83(7-8):1102–1116.
- 493 [30] Moreno, A., Poch, D., Bonafonte, A., Lleida, E., Llisterri, J., Marino, J. B., and Nadeu,
494 C. (1993). Albayzín speech database: Design of the phonetic corpus. In *Eurospeech*, pages
495 175–178, Berlin, Germany.
- 496 [31] Naylor, G. (2016). Theoretical issues of validity in the measurement of aided speech reception
497 threshold in noise for comparing nonlinear hearing aid systems. *J. Am. Acad. Audiology*,
498 27:504–514.
- 499 [32] Nejime, Y. and Moore, B. C. J. (1998). Evaluation of the effect of speech-rate slowing on
500 speech intelligibility in noise using a simulation of cochlear hearing loss. *J. Acoust. Soc. Am.*,
501 103(1):572–576.
- 502 [33] Oxenham, A. J. and Simonson, A. M. (2009). Masking release for low-and high-pass-filtered
503 speech in the presence of noise and single-talker interference. *J. Acoust. Soc. Am.*, 125:457–468.
- 504 [34] Patterson, R. D., Holdsworth, J., Nimmo-Smith, I., and Rice, P. (1988). SVOS Final Report:
505 The Auditory Filterbank. Technical report 2341. MRC Applied Psychology Unit.
- 506 [35] Picheny, M. A., Durlach, N. I., and Braida, L. D. (1985). Speaking clearly for the hard of
507 hearing. I: Intelligibility differences between clear and conversational speech. *J. Speech Hear.*
508 *Res.*, 28:96–103.
- 509 [36] Pisoni, D. B., Bernacki, R. H., Nusbaum, H. C., and Yuchtman, M. (1985). Some acoustic-
510 phonetic correlates of speech produced in noise. In *ICASSP*, pages 1581–1584, Tampa, Florida.
- 511 [37] Pittman, A. L. and Wiley, T. L. (2001). Recognition of speech produced in noise. *J. Speech*
512 *Lang. Hear. Res.*, 44(3):487–496.
- 513 [38] Rothauser, E. H., Chapman, W. D., Guttman, N., Hecker, M. H. L., Nordby, K. S., Silbiger,

- 514 H. R., Urbanek, G. E., Weistock, M., McGee, V. E., Pahl, U. P., and Voiers, W. D. (1969).
515 IEEE Recommended practice for speech quality measurements. *IEEE Trans. Audio Acoust.*,
516 pages 225–246.
- 517 [39] Sankowska, J., García Lecumberri, M. L., and Cooke, M. (2011). Interaction of intrinsic
518 vowel and consonant durational correlates with foreigner directed speech. *Poznań Studies in
519 Contemporary Linguistics*, 47:109–119.
- 520 [40] Sauert, B. and Vary, P. (2006). Near end listening enhancement: Speech intelligibility im-
521 provement in noisy environments. In *Proc. ICASSP*, pages 493–496, Toulouse, France.
- 522 [41] Schepker, H., Rennie, J., and Doclo, S. (2013). Improving speech intelligibility in noise
523 by sii-dependent preprocessing using frequency-dependent amplification and dynamic range
524 compression. In *Proc. Interspeech*, pages 3577–3581.
- 525 [42] Skowronski, M. D. and Harris, J. G. (2006). Applied principles of clear and Lombard speech
526 for automated intelligibility enhancement in noisy environments. *Speech Communication*,
527 48(5):549–558.
- 528 [43] Song, J. Y., Demuth, K., and Morgan, J. (2010). Effects of the acoustic properties of infant-
529 directed speech on infant word recognition. *J. Acoust. Soc. Am.*, 128(1):389–400.
- 530 [44] Stilp, C. and Kluender, K. (2010). Cochlea-scaled entropy, not consonants, vowels, or time,
531 best predicts speech intelligibility. *P. Natl. Acad. Sci. USA*, 107(27):12387–12392.
- 532 [45] Summers, W. V., Pisoni, D. B., Bernacki, R. H., Pedlow, R. I., and Stokes, M. A. (1988).
533 Effects of noise on speech production: Acoustic and perceptual analyses. *J. Acoust. Soc. Am.*,
534 84(3):917–928.
- 535 [46] Taal, C. H., Jensen, J., and Leijon, A. (2013). On optimal linear filtering of speech for near-end
536 listening enhancement. *IEEE Signal Proc. Let.*, 20(3):225–228.
- 537 [47] Tang, Y. and Cooke, M. (2012). Optimised spectral weightings for noise-dependent speech
538 intelligibility enhancement. In *Proc. Interspeech*, pages 955–958, Portland, USA.
- 539 [48] Uchanski, R. M. (2005). Clear speech. In Pisoni, D. B. and Remez, R. E., editors, *The
540 Handbook of Speech Perception*, pages 207–235. Blackwell, Oxford, UK.
- 541 [49] Valentini-Botinhao, C., Yamagishi, J., and King, S. (2012). Mel cepstral coefficient modifi-
542 cation based on the Glimpse Proportion measure for improving the intelligibility of HMM-
543 generated synthetic speech in noise. In *Proc. Interspeech*, pages 631–634, Portland, USA.
- 544 [50] Yoo, S. D., Boston, J. R., El-Jaroudi, A., Li, C.-C., Durrant, J. D., Kovacyk, K., and Shaiman,

- 545 S. (2007). Speech signal modification to increase intelligibility in noisy environments. *J.*
546 *Acoust. Soc. Am.*, 122(2):1138–1149.
- 547 [51] Zorila, T., Kandia, V., and Stylianou, Y. (2012). Speech-in-noise intelligibility improvement
548 based on spectral shaping and dynamic range compression. In *Proc. Interspeech*, pages 635–
549 638.
- 550 [52] Zorila, T.-C. and Stylianou, Y. (2015). A fast algorithm for improved intelligibility of speech-
551 in-noise based on frequency and time domain energy reallocation. In *Proc. Interspeech*, pages
552 60–64.