

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

Martin Cooke^{a)}

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

Vincent Aubanel

University of Grenoble Alpes, Centre National de la Recherche Scientifique, GIPSA-lab, Grenoble, France

(Received 17 January 2017; revised 5 May 2017; accepted 8 May 2017; published online 5 June 2017)

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. © 2017 Acoustical Society of America.

[<http://dx.doi.org/10.1121/1.4983826>]

[DB]

Pages: 4126–4135

I. INTRODUCTION

In speech communication scenarios involving the output of natural or synthetic speech, the likelihood of correct message reception in noisy environments can be improved by modifying the speech signal prior to output (e.g., Skowronski and Harris, 2006; Sauert and Vary, 2006; Yoo *et al.*, 2007; Brouckxon *et al.*, 2008; Valentini-Botinhao *et al.*, 2012; Taal *et al.*, 2013; Zorila and Stylianou, 2015; Jokinen *et al.*, 2016). Such approaches are highly effective: A recent evaluation of 18 algorithms demonstrated gains equivalent to increasing signal-to-noise ratio (SNR) by 5.1 dB for natural speech and by 5.6 dB for synthetic speech (The Hurricane Challenge; Cooke *et al.*, 2013).

Most speech modification techniques operate by reallocating energy in time and frequency under a constant input-output root-mean-square (RMS) energy constraint. Energy reallocation aims to enhance intelligibility by manipulating the spectro-temporal pattern of local SNR, enabling weaker regions to rise above the masker at the expense of portions of the signal whose local SNR is already sufficiently high. Different approaches have variously transferred energy from voiced to voiceless regions (Skowronski and Harris, 2006), boosted some regions of the spectrum at the expense of others (Tang and Cooke, 2012), enhanced formants (Brouckxon *et al.*, 2008), increased the amplitude modulation

depth of the mid-frequencies (Koutsogiannaki and Stylianou, 2016), or employed dynamic range compression (Blesser, 1969), which has the effect of transferring energy from intense to weaker temporal epochs (Zorila *et al.*, 2012; Schepker *et al.*, 2013). Cooke *et al.* (2014a) provides a review of human and algorithmic speech modifications.

An alternative to spectro-temporal energy reallocation is the modification of segment or sub-segmental durations. Altered speaking styles such as Lombard speech (e.g., Pisoni *et al.*, 1985; Summers *et al.*, 1988; Junqua, 1993), clear speech (Picheny *et al.*, 1985; Uchanski, 2005), speech directed at infants (e.g., Grieser and Kuhl, 1988), and speech produced at a distance (Fux *et al.*, 2012) exhibit durational changes, usually resulting in slower speech, both overall and at the level of individual speech segments. Many of these forms of speech have been found to be more intelligible than unmodified plain speech (Dreher and O'Neill, 1957; Picheny *et al.*, 1985; Pittman and Wiley, 2001; Lu and Cooke, 2008; Song *et al.*, 2010). While acoustic changes to features such as intensity, spectral tilt, and prosody are present in altered speech styles and may play a role in increased intelligibility, it is natural to consider whether durational changes contribute to the improvement. Durational manipulations can be expected to be particularly effective in the presence of fluctuating maskers where the opportunity arises to shift phonetic information in time to regions where the masking source is less intense.

Approaches employing durational modifications are less common than those exploiting spectro-temporal energy

^{a)}Also at Ikerbasque (Basque Science Foundation). Electronic mail: m.cooke@ikerbasque.org

reallocation. Of the 14 natural speech modification approaches submitted to the aforementioned Hurricane Challenge (Cooke *et al.*, 2013), only 3 involved significant durational changes. Tellingly, two of these three produced the largest enhancements in the fluctuating masker condition of the Challenge at moderate and adverse SNRs. Indeed, gains of up to 4.4 dB resulted from an approach (GCReTime; Aubanel and Cooke, 2013) that modified durational information only, indicating that alterations to segment durations alone can be a valuable strategy for maskers containing low-frequency temporal modulations.

However, the basis for the intelligibility enhancements produced by durational changes is currently unclear. It is possible that listeners are able to take advantage of the reduced information rate of slower speech rate rather than the intended energetic masking release produced by shifting information in time. Evidence for intelligibility benefits of speech rate slowing is mixed. While studies by Adams and colleagues (Adams and Moore, 2009; Adams *et al.*, 2012) have demonstrated intelligibility increases for slow speech in masking noise, no such effect was observed under conditions of simulated hearing loss by Nejime and Moore (1998) nor when linear and nonlinear durational changes observed in Lombard speech were mapped on to plain speech (Cooke *et al.*, 2014b). Intriguingly, while the latter studies used stationary maskers, the sentence material used by Adams and Moore (2009) and Adams *et al.* (2012) was mixed with four-talker babble, leading to the possibility that the temporal modulation characteristics of the masker played a role in the different outcomes.

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

The current study also addressed the issue of whether the intelligibility of durationally modified speech is affected by the properties of the masker. Utterances were presented in stationary noise and two forms of fluctuating noise: competing speech (CS) and speech-modulated noise (SMN) with temporal envelope fluctuations matching those of CS.

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

A. Durational modifications

Listeners heard sentences that were either unmodified (PLAIN), linearly stretched (ELONGATED) or nonlinearly modified (RETIMED). All durational modifications were carried out using the WSOLA algorithm (Demol *et al.*, 2005) via a sequence of time-scale factors. In the elongation condition, a constant time-scale factor was used, while in the retiming case the time-scale factor sequence was derived from the

GCReTime algorithm described in Aubanel and Cooke (2013) and summarised below.

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

The GCReTime local distance function $D(i,j)$ is defined on a grid of points i,j corresponding to the i th frame of the target speech signal s and the j th frame of the masker m . The local distance function for all possible pairs (i,j) is a matrix, shown as a grayscale image in Fig. 1, where darker regions depict higher values of the function. The local distance function is composed of two components quantifying (1) the masked audibility of the speech signal in frame i in the

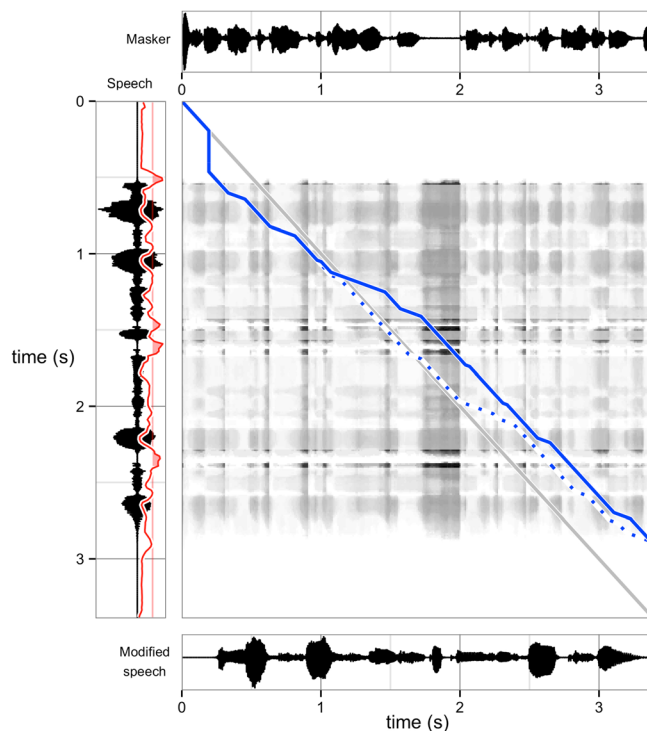


FIG. 1. (Color online) An illustration of sentence retiming in the face of a CS masker. The grayscale image depicts the value of the local cost function [Eq. (A1)] for all possible pairs of frames of the target and masker. The solid line shows the minimum cost retiming path using the glimpse proportion (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 the Appendix).

presence of the masker at frame j , and (2) the informativeness of the speech signal in the vicinity of frame i . The first of these components is operationalised using the glimpse proportion (GP; [Cooke, 2006](#)), while the second makes use of cochlear-scaled entropy (CSE; [Stilp and Kluender, 2010](#)); together these components are reflected in the name “GCRetime.” CSE captures localised spectral change and has been shown to predict intelligibility better than consonants, vowels, or consonant-vowel/vowel-consonant transitions when tested using a noise-replacement paradigm ([Stilp and Kluender, 2010](#)). Taken together, these two components ensure that the distance function takes on high values when the speech signal is not masked and when it is undergoing a period of rapid change. For example, the dark vertical band in the period immediately preceding the 2 s point in the masker is due to the low level of the masker in that interval, and the darker horizontal strips within this band correspond to those portions of the target speech with a high CSE value. The path that maximises the global distance passes through some of these regions, effectively ensuring that potentially high-information-value transients in the target speech are shifted in time to regions where the masker is less intense.

The [Appendix](#) describes the computation of the GCRetime local distance function D in more detail.

B. Speech and masker materials

Utterances were drawn from the phonemically balanced Sharvard corpus ([Aubanel et al., 2014](#)), which consists of Spanish sentences designed to be equivalent in difficulty to the Harvard sentences ([Rothausen et al., 1969](#)). Each Sharvard sentence contains five keywords used for scoring; an example (keywords underlined) is “Llene el frasco de crystal con cola densa” (“Fill the glass flask with thick glue”). Spectrograms of this utterance in each of the three styles PLAIN, ELONGATED, and RETIMED are shown in Fig. 2. Renditions of the first 243 Sharvard utterances read by a native Spanish male talker were used in experiment 1; this number includes sentences used as practice items.

Maskers were constructed using speech material from a native Spanish female talker who read sentences from the Albayzin corpus ([Moreno et al., 1993](#)). Inter-sentence pauses were removed and sentences concatenated to produce a signal of 13.83 min duration, sufficient to ensure that no masker fragment was repeated in any speech-plus-masker mixture. Successive non-overlapping fragments from this signal were used for the CS masking condition. A speech-shaped noise (SSN) masker was constructed by passing white noise through a filter with a long-term spectrum matching that of the female talker. Each CS fragment had a matched SMN signal formed by multiplying the short-term temporal envelope of the CS fragment with a portion of the SSN signal selected at random. All speech and masker materials were sampled at 16 kHz.

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

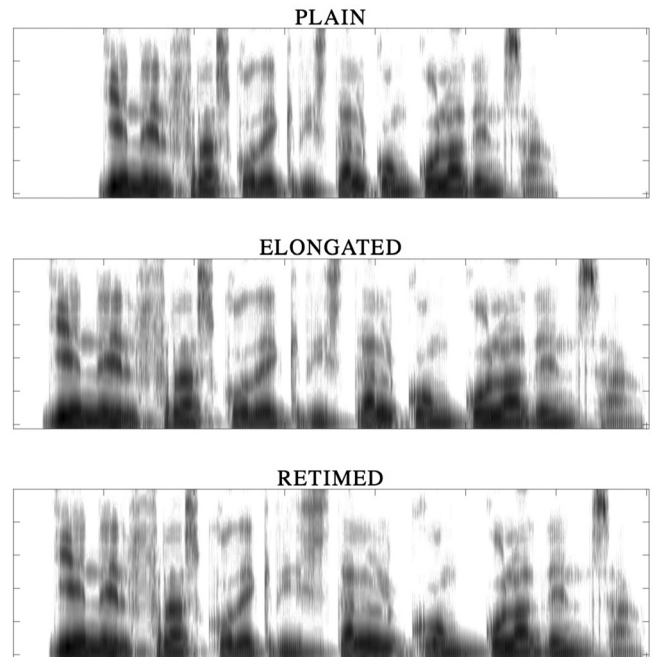


FIG. 2. Example spectrograms of the utterance “Llene el frasco de cristal con cola densa” in each of the three durational modification conditions.

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

C. Participants

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

D. Procedure

Listeners heard a total of 234 utterances made up of 26 sentences in each of the 9 conditions resulting from the combination of the PLAIN, ELONGATED, and RETIMED manipulations with the 3 maskers (CS, SMN, and SSN). The SNR for the SSN masking conditions was set to -6.5 dB, a value which led to a 50% keyword score for the male talker in [Aubanel et al. \(2014\)](#). Since CS is a less effective masker than SSN when presented at the same SNR, the SNR for the CS masker was set following pilot tests to -17 dB, while similar tests indicated the need for an intermediate SNR for the SMN case of -12 dB.

To avoid sentence subset effects due to possible differing intrinsic intelligibilities of the speech material, the three speech processing conditions (PLAIN, ELONGATED, RETIMED)

were applied to the complete set of 234 utterances. Listeners were assigned to subsets of sentences in such a way as to ensure that each sentence in each processing condition was heard the same number of times across listeners, and that each listener heard each sentence exactly once. Speech-plus-noise stimuli were blocked by masker type; within each block listeners heard equal numbers of sentences from each of the three processing conditions in a randomised order. Immediately prior to each block of 78 sentences, listeners responded to 3 unscored practice stimuli designed to familiarise them with the type of masker. Practice sentences did not occur elsewhere in the main experiment. Presentation order of the three blocks was balanced across listeners.

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

E. Postprocessing

Listener responses were scored automatically based on the number of keywords identified correctly in each sentence. Vowel stress marks were removed prior to scoring, so that, for example, both “mas” and “más” were considered to be correct responses for the word “más.” A total of 130 keywords were scored (26 sentences \times 5 keywords per sentence) in each of the 9 conditions. Scores were expressed as percentages of keywords identified correctly in each condition. Since none of the scores lay outside the range 23%–83%, raw (untransformed) percentages were used in subsequent statistical analyses.

F. Results

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

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

Retimed speech produced a larger spread of differences over the PLAIN baseline across the three maskers. In stationary noise, retiming was highly detrimental to intelligibility, producing a loss of 14.9 p.p. compared to unmodified speech. For the modulated noise masker (SMN) the gain of 10.3 p.p. was similar to that seen for the ELONGATED condition. However, with a gain of 16.3 p.p., retimed utterances

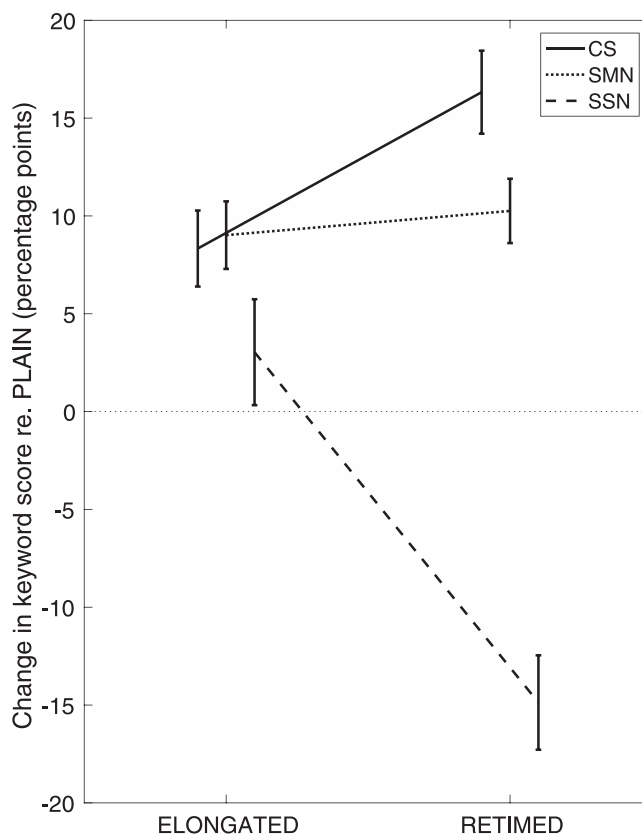


FIG. 3. Changes in mean keywords correct relative to PLAIN speech for ELONGATED and RETIMED utterances in the presence of CS, SMN, and SSN maskers. Error bars here and in Fig. 5 represent ± 1 standard error.

were substantially more intelligible than their elongated counterparts in the CS masker.

An analysis of variance (ANOVA) of changes-over-baseline scores with within-subjects factors of modification method (ELONGATED, RETIMED) and masker (CS, SMN, SSN) demonstrated a clear interaction in the effect of modifications and maskers [$F(2,34) = 37.7, p < 0.001, \eta^2 = 0.28$, mean square error (MSE) = 43.1], with significant main effects of both modification type [$F(1,17) = 10.8, p < 0.01, \eta^2 = 0.03$, MSE = 20.8], and masker [$F(2,34) = 21.2, p < 0.001, \eta^2 = 0.46$, MSE = 164.4]. Based on a Fisher’s least significant difference (LSD) of 4.4 p.p., gains for ELONGATED speech in the two modulated maskers were equivalent, while ELONGATED speech was statistically equivalent to the PLAIN baseline for the SSN masker. Changes in keyword scores for RETIMED speech were significantly different in the three masking conditions.

G. Interim discussion

A strategy of retiming speech information by shifting the waveform nonlinearly in time to attenuate the effect of intense masker epochs has previously been shown to produce substantial intelligibility gains in the presence of fluctuating maskers (Aubanel and Cooke, 2013). Experiment 1 confirms the effectiveness of algorithmic speech retiming, and extends this finding to speech material in a different language: the 16.3 p.p. gain produced for Spanish sentences in the RETIMED condition of the current study in the CS masker condition at an SNR of -17 dB is consistent with the improvements of

16 and 18 p.p. observed in Aubanel and Cooke (2013) for English sentences at SNRs of -14 and -21 dB in the equivalent masking condition of that study.

Elongation of speech had a negligible impact on intelligibility for stationary maskers, suggesting that a slower speech rate in itself is not responsible for the gains observed when speech is retimed. In contrast, for fluctuating maskers, elongation led to a clear increase in intelligibility. This finding goes some way to explaining the discrepancies among earlier studies on the effectiveness of a slower speech rate in noise. As noted in the Introduction, while Nejime and Moore (1998) and Cooke *et al.* (2014b) failed to find an intelligibility benefit of slower speech when presented in a stationary masker, Adams *et al.* (2012) reported a beneficial effect of a slower speech rate in four-talker babble, a type of masker that shows a greater temporal modulation depth than that of a purely stationary noise. The issue of how a fluctuating masker might promote intelligibility increases for elongated speech is addressed in Sec. IV, General Discussion.

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

III. EXPERIMENT 2: ROLE OF THE RETIMING MASKER

A. Listeners

A new cohort of 21 normal hearing paid native Spanish speakers (16 female) with a mean age of 20.0 years (s.d. = 1.5) and the same profile as the participants of experiment 1 took part in experiment 2. Results from one participant who treated the CS masker as the target in a number of conditions were excluded.

B. Materials and methods

Raw speech materials, maskers, and SNRs were the same as those used in experiment 1. The ELONGATED condition was not tested. Instead, listeners heard utterances in unmodified PLAIN form and in two distinct retiming conditions. In one retiming condition (CS RETIMED) the durational

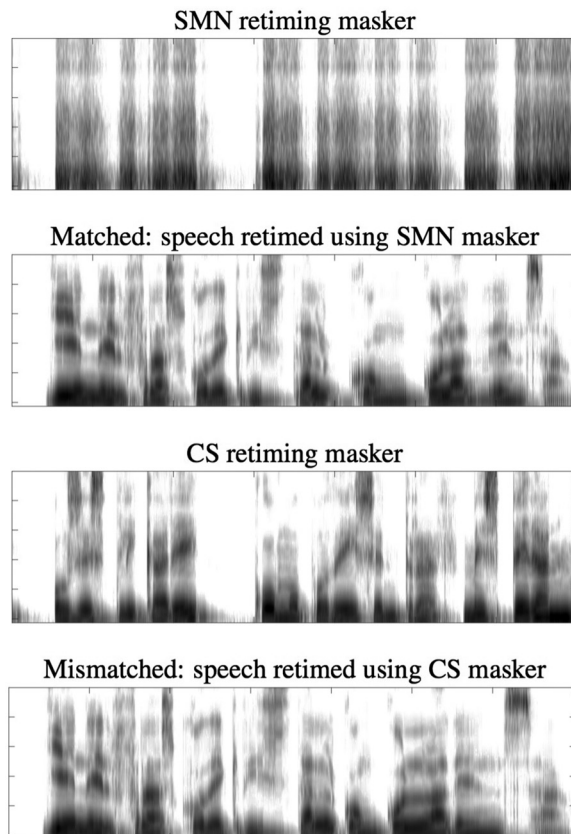


FIG. 4. Matched and mismatched retiming for the utterance shown in Fig. 2. The top panel shows the SMN masker used to produce the retiming path, which results in the utterance shown in the second panel (SMN RETIMED). The third panel shows the CS masker used for retiming which results in the utterance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

modifications were based on GCReTime using a CS masker, while in the other retiming condition (SMN RETIMED) the modifications result from the counterpart SMN masker. In this way, listeners heard sentences retimed by a matched or unmatched masker. The nine experimental conditions (3 modifications \times 3 maskers) were presented to listeners using the blocking and balancing procedure of experiment 1 as described in Sec. IID.

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

C. Results

Mean keywords correct scores for the PLAIN speech condition were 49.6%, 42.5%, and 55.8% for the CS, SMN, and SSN presentation maskers, respectively. Figure 5 plots changes in scores over the PLAIN baseline for sentences retimed using the CS and SMN maskers for each of the three presentation maskers. Changes over baseline for the CS RETIMED conditions were similar to those observed in the equivalent conditions of experiment 1 (CS, 16.3 vs 16.5 p.p.;

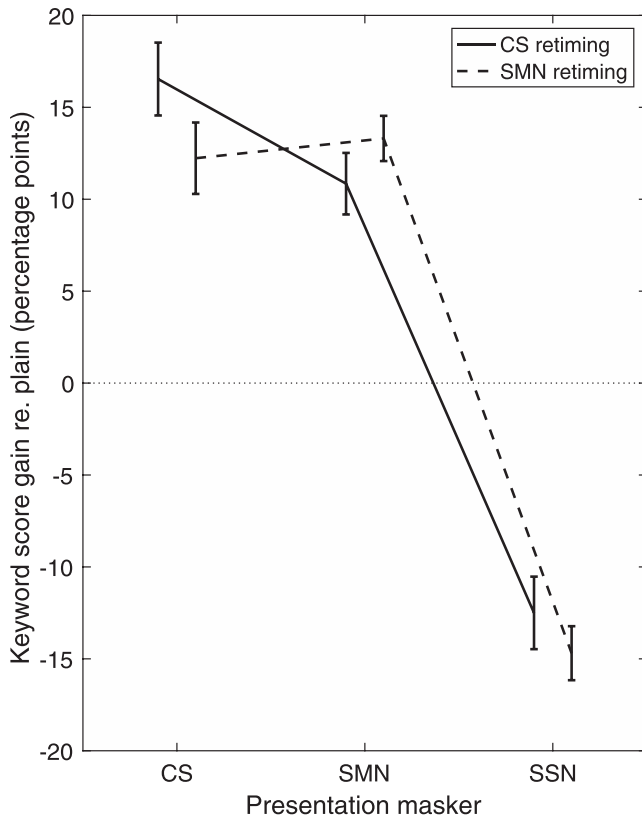


FIG. 5. Changes in keywords correct scores relative to unmodified plain speech for utterances retimed using the CS or SMN maskers for the three presentation maskers.

SMN 10.2 vs 10.8; SSN -14.9 vs -12.5), confirming the findings of the first experiment with a different listener cohort.

A within-subjects ANOVA with factors of retiming masker and presentation masker for the two fluctuating masking conditions (CS and SMN) indicated no main effect of either factor, but revealed a significant interaction between the two factors [$F(1, 19) = 6.96, p < 0.05, \eta^2 = 0.048, MSE = 32.9$]. *Post hoc* tests based on a Fisher's LSD of 3.80 p.p. indicate that CS-based retiming was more effective in a matched CS masker (16.5 p.p.) than in a mismatched SMN masker (10.9 p.p.). However, there was no benefit of matched masker type for SMN-based retiming, with similar gains of 13.3 and 12.2 p.p. in the matched and unmatched conditions, respectively. Critically, CS RETIMED speech led to higher gains than SMN RETIMED speech when presented in a CS masker, suggesting that the specific details of the retiming path are important. Gains in the two matched conditions (i.e., CS RETIMED in CS masker and SMN RETIMED in SMN masker) did not differ statistically, the difference of 3.2 p.p. falling short of the critical LSD value.

D. Interim discussion

Experiment 2 demonstrated that the benefits of retiming are affected by the relationship between the retiming masker and the presentation masker in the case of the CS masker but not for the SMN masker. This outcome suggests that there is a limit to the benefits of retiming for a temporally modulated

noise masker, while for CS there may be both temporal and spectral opportunities which are taken into account by the energetic masking model underlying the GP calculation.

IV. GENERAL DISCUSSION

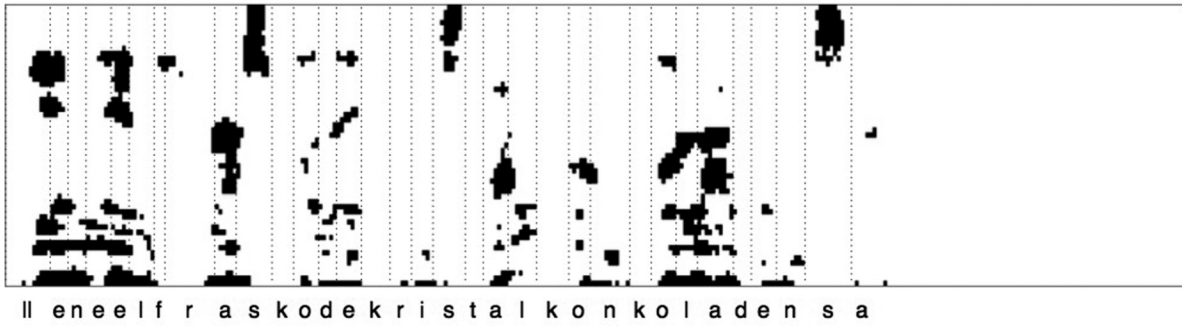
Experiment 1 addressed the primary research question of the current study by measuring the extent to which intelligibility gains are present for speech that is linearly elongated to generate the same average speech rate as that produced by retiming. The absence of a benefit of elongated speech in the presence of a stationary speech-shaped masker appears to rule out reduced speech rate *per se* as a contributory factor.

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

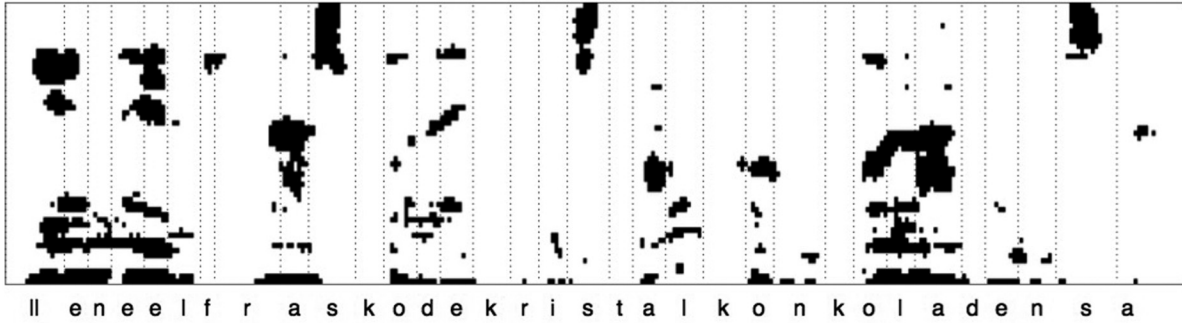
There are several ways in which a temporally fluctuating masker might promote intelligibility gains for slowed speech while a stationary masker does not. One possibility is that the regions of elongated speech which escape masking by a stationary noise provide no new phonetic information. The upper two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

An alternative explanation for the observed gains in fluctuating maskers arises from the possibility that listeners are better able to separate target and background speech due to speech rate differences between the target and masker, overcoming a potential source of informational masking. This notion is supported by a study by Gordon-Salant and Fitzgibbons (2004) in which a cohort of young normal hearing listeners recognised more words in time-compressed sentences when the compression ratio did not match that of 12-talker background babble. The ELONGATED condition of the current study does indeed lead to an increase in speech rate differences between the target and background: The PLAIN speech was articulated at a rate of 5.8 vowels/s, comparable

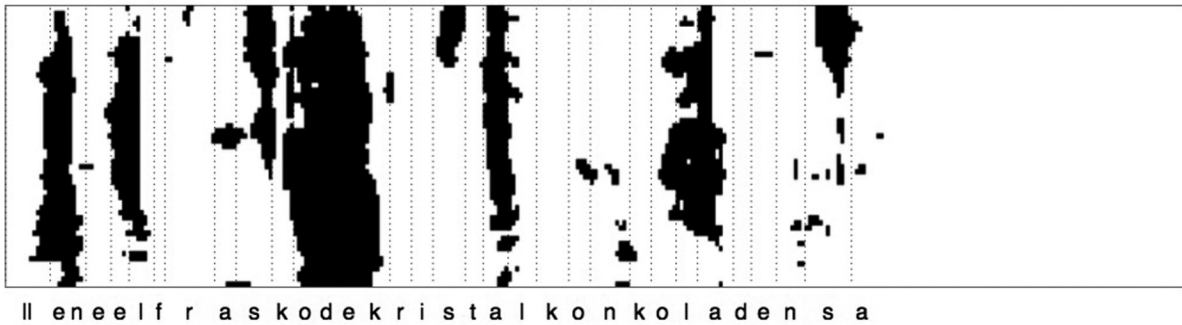
PLAIN in SSN



ELONGATED in SSN



PLAIN in CS



ELONGATED in CS

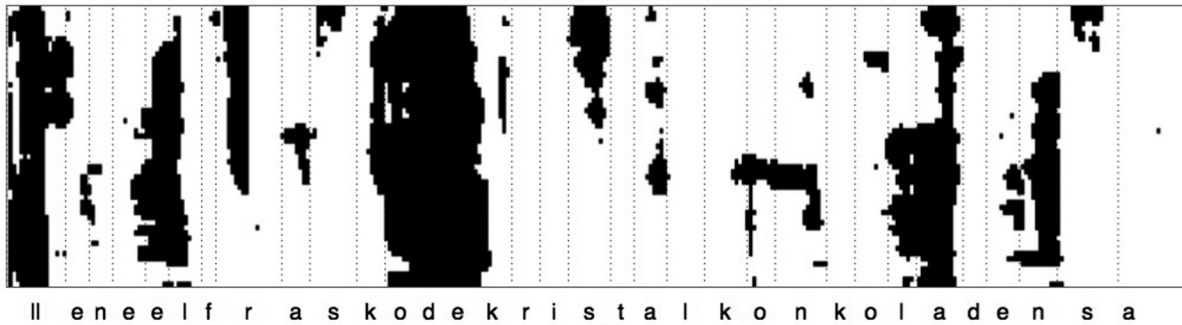


FIG. 6. Regions of a target speech utterance of the phrase “Llene el frasco de cristal con cola densa” which are deemed to escape energetic masking according to a glimpsing model (Cooke, 2006), for PLAIN and ELONGATED speech in the presence of a stationary masker (top two panels) and a fluctuating masker (lower two panels). A broad phoneme-level transcription is provided in each case.

to the 6.6 vowels/s of the CS masker. Speech rate slowing in the ELONGATED condition reduced the average speech rate to 4.4 vowels/s, increasing the target-background speech rate difference. Earlier studies with listeners (Miller and Licklider, 1950) or models (Bronkhorst *et al.*, 1993) have demonstrated

sensitivity to differences in speech interruption rate and speech modulation rate. Further studies controlling for speech rate differences are needed to rule out their possible contribution to the observed gains in the current study. However, since gains were also observed for the non-speech masker, it

appears necessary to invoke a more generalised notion of temporal modulation rate differences that go beyond speech-on-speech informational masking, which in this case is likely to have a relatively small effect since the gender of the target and masking talkers differed (Brungart *et al.*, 2001). A related possibility is that the presence of modulation in the masker itself imposes a cognitive burden on listeners: A slower rate of information transmission via a lower speech rate may be beneficial in reducing listening effort.

A striking and consistent outcome observed in both experiments was the substantial loss of intelligibility (amounting to 13–15 p.p.) that occurred when speech was nonlinearly retimed in stationary noise, contrasting with no loss for linearly elongated speech in the presence of the stationary masker. In the GCRetime algorithm, retiming takes no account of anything other than the temporal relationship between speech and masker (the glimpsing component) and speech dynamics (the CSE component). Consequently, properties such as segment duration and the local speech rate of unmodified speech are not preserved by the algorithm. Such distortions are likely both to confound listeners' expectations of when salient information is going to occur and to diminish the effectiveness of contextual cues that depend on relative durations. For Spanish, changes in relative segment durations induced by GCRetime may have interfered with phonological cues (e.g., Mendoza *et al.*, 2003) or syllabification (e.g., Hualde and Chitoran, 2003). Intriguingly, such a tradeoff between modifications that overcome masking and those which preserve phonological integrity is also seen in naturally produced speech in noise. Sankowska *et al.* (2011) measured durational (vowel lengthening) cues to the voicing distinction in English plosives in plain and Lombard speech, finding a reduced contrast in the latter case. The benefits of retiming in fluctuating maskers presumably reflect a net effect of masking release and durational distortion, suggesting that even larger benefits in noise are realisable if the phonological impact of durational modifications can be minimised.

Distortions to the target speech might also have contributed to the observed differences in effectiveness of retiming in the presence of CS maskers and temporally modulated noise, even when the influence of a matching or mismatching retiming masker was controlled for, as in experiment 2. In a non-informational masker such as modulated noise, listeners' attentional focus is presumably directed to the target speech alone, and any departures from expected phonological forms may be noticeable, and potentially lead to the consideration of additional competitor words. In contrast, when the masker itself contains speech, it is conceivable that the cognitive burden imposed by foreground-background separation precludes a more detailed analysis of the target signal, or a mis-attribution of retiming-based distortions to the CS signal. Another possibility is that any gains due to retiming outweigh losses due to mistiming of phonological features.

The outcome of the current study points to the potential of durational changes as a mechanism for improving intelligibility in noise, but also highlights the need to take the temporal properties of the masker into account, given the deficits resulting from the retiming method in the presence of stationary noise. The finding that gains are possible

merely by elongating the speech signal in fluctuating maskers suggests that speech rate slowing could be a component of a simple practical strategy for boosting intelligibility. As noted in the Introduction, modified duration is not by any means the sole manifestation of natural "altered" speaking styles, and spectral factors in particular are known to have a sizeable influence on intelligibility (Cooke *et al.*, 2014b). Spectral and durational changes are orthogonal to a large extent, e.g., changes to properties such as spectral tilt can be imposed independently of durational changes.

As is typically the case when targeting the 50% correct response rate with a normal-hearing adult population, all testing was done at negative SNRs. Further work is needed to measure the efficacy of a slower speech rate at more realistic SNRs (Naylor, 2016), as such environments have also been found to induce slower rate of speech in talkers (Aubanel *et al.*, 2011). The benefits observed in the current study of nonlinear retiming at negative SNRs may be reduced at higher SNRs; lower than expected benefits for a fluctuating masker advantage in comparison to stationary noise have consistently been observed at positive SNRs (Bernstein and Grant, 2009; Oxenham and Simonson, 2009; Freyman *et al.*, 2008).

V. CONCLUSIONS

- (i) Reductions in speech rate resulting from linear elongation of the speech signal did not lead to intelligibility increases (nor did they disrupt intelligibility) for sentences in the presence of stationary speech-shaped maskers, suggesting that intelligibility gains seen in durationally modified speech were not due to the reduction in information rate that accompanies slower speech.
- (ii) However, identical elongations produced significant intelligibility increases in fluctuating maskers. One explanation is that while elongation in stationary maskers produces no new speech information, the altered pattern of glimpses in fluctuating maskers leads to the unmasking of new phonetic cues. An alternative is that slower speech enables listeners to separate target speech from the background due to greater differences in speech rate, or reduces the cognitive burden of processing speech in a modulated background. Further studies are needed to distinguish these possibilities.
- (iii) Nonlinear durational modifications designed to reduce energetic masking of speech information led to larger intelligibility gains in CS maskers than those produced by linear elongation in spite of the distortion of phonetic integrity indicated by the reduced intelligibility of the same modifications in stationary maskers.

ACKNOWLEDGMENTS

This work was partially funded by the "Listening Talker" project, supported by the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European

Commission under FET-Open Grant No. 256230. V.A. also acknowledges support from the FP7 FET Project ‘‘Speech Unit(e)s,’’ Grant No. 339152.

APPENDIX: COMPUTATION OF THE GCReTime LOCAL DISTANCE FUNCTION

The GCReTime local distance function [Eq. (A1)] is defined for each pair of time frames i of the speech signal and j of the masker as the product of two terms: (i) GP, $GP(i, j)$, the proportion of the speech signal in frame i glimpsed in the presence of the masker in frame j [Eq. (A2)]; and (ii) $W_{\text{CSE}}(i)$, a weighting term based on the CSE of the speech signal in frame i [Eq. (A3)]

$$D(i, j) = GP(i, j)W_{\text{CSE}}(i). \quad (\text{A1})$$

1. GP

The GP is intended to reflect the local audibility of speech in noise and is defined as the percentage of spectral regions where the modelled auditory excitation pattern for the target speech exceeds that of the masker

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

where F is the number of frequency channels, S_f and M_f denote the excitation patterns of speech and masker in frequency channel f , respectively, and $\mathcal{H}(\cdot)$ is the Heaviside unit step function counting the number of channels where the speech exceeds the masker. Excitation patterns are derived via a gammatone filterbank (Patterson *et al.*, 1988) using an implementation introduced by Cooke (1993). The Hilbert envelope of each gammatone filter output is computed and smoothed by a leaky integrator with an 8 ms time constant (Moore *et al.*, 1988), downsampled and log-compressed. Here, the gammatone filterbank contained $F = 32$ frequency channels spaced equally on an equivalent rectangular bandwidth-rate scale between 50 Hz and 7500 Hz.

2. CSE

In the current study we use the concept of CSE (Stilp and Kluender, 2010) to identify spectral regions which are changing most rapidly in order to give them greater weight in the computation of the local distance function. CSE is implemented as a locally averaged measure of spectral change across time based on excitation patterns of the target speech signal

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

where

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

and λ is the number of frames over which the CSE is computed. Following Stilp and Kluender (2010), $\lambda = 5$, equivalent to 80 ms for the 16 ms time frames used here.

The CSE-based weighting is defined as

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

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

- Adams, E. M., Gordon-Hickey, S., Morlas, H., and Moore, R. (2012). ‘‘Effect of rate-alteration on speech perception in noise in older adults with normal hearing and hearing impairment,’’ *Am. J. Audiol.* **21**(1), 22–32.
- Adams, E. M., and Moore, R. E. (2009). ‘‘Effects of speech rate, background noise, and simulated hearing loss on speech rate judgment and speech intelligibility in young listeners,’’ *J. Am. Acad. Audiol.* **20**, 28–39.
- Aubanel, V., and Cooke, M. (2013). ‘‘Information-preserving temporal reallocation of speech in the presence of fluctuating maskers,’’ in *Proc. Interspeech*, Lyon, France, pp. 3592–3596.
- Aubanel, V., Cooke, M., Villegas, J., and Garca Lecumberri, M. L. (2011). ‘‘Conversing in the presence of a competing conversation: Effects on speech production,’’ in *Proc. of Interspeech*, Florence, Italy, pp. 2833–2836.
- Aubanel, V., Garca Lecumberri, M. L., and Cooke, M. (2014). ‘‘The Sharvard Corpus: A phonemically-balanced Spanish sentence resource for audiology,’’ *Int. J. Audiol.* **53**, 633–638.
- Bernstein, J. G., and Grant, K. W. (2009). ‘‘Auditory and auditory-visual intelligibility of speech in fluctuating maskers for normal-hearing and hearing-impaired listeners,’’ *J. Acoust. Soc. Am.* **125**, 3358–3372.
- Blessner, B. A. (1969). ‘‘Audio dynamic range compression for minimum perceived distortion,’’ *IEEE Trans. Audio Electroacoust.* **17**(1), 22–32.
- Bronkhorst, A. W., Bosman, A. J., and Smoorenburg, G. F. (1993). ‘‘A model for context effects in speech recognition,’’ *J. Acoust. Soc. Am.* **93**, 499–509.
- Brouckxon, H., Verhelst, W., and Schuymer, B. D. (2008). ‘‘Time and frequency dependent amplification for speech intelligibility enhancement in noisy environments,’’ in *Proc. Interspeech*, Vol. 9, pp. 557–560.
- Brungart, D. S., Simpson, B. D., Ericson, M. A., and Scott, K. R. (2001). ‘‘Informational and energetic masking effects in the perception of multiple simultaneous talkers,’’ *J. Acoust. Soc. Am.* **110**(5), 2527–2538.
- Cooke, M. (1993). *Modelling Auditory Processing and Organisation* (Cambridge University Press, Cambridge).
- Cooke, M. (2006). ‘‘A glimpsing model of speech perception in noise,’’ *J. Acoust. Soc. Am.* **119**(3), 1562–1573.
- 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., Mayo, C., and Valentini-Botinhao, C. (2013). ‘‘Intelligibility-enhancing speech modifications: The Hurricane Challenge,’’ in *Proc. Interspeech*, pp. 3552–3556.
- Cooke, M., Mayo, C., and Villegas, J. (2014b). ‘‘The contribution of durational and spectral changes to the Lombard speech intelligibility benefit,’’ *J. Acoust. Soc. Am.* **135**(2), 874–883.
- Demol, M., Verhelst, W., Struyve, K., and Verhoeve, P. (2005). ‘‘Efficient non-uniform time-scaling of speech with WSOLA,’’ in *Int. Conf. on Speech and Computers (SPECOM)*, pp. 163–166.
- Dreher, J. J., and O’Neill, J. J. (1957). ‘‘Effects of ambient noise on speaker intelligibility for words and phrases,’’ *J. Acoust. Soc. Am.* **29**(12), 1320–1323.
- Freyman, R. L., Balakrishnan, U., and Helfer, K. S. (2008). ‘‘Spatial release from masking with noise-vocoded speech,’’ *J. Acoust. Soc. Am.* **124**, 1627–1637.
- Fux, T., Feng, G., and Zimpfer, V. (2012). ‘‘Natural-to-shouted voice transformation for distance cues of monosyllabic consonant-vowel-consonant words,’’ *Acta Acust. Acust.* **98**(5), 839–843.
- Gordon-Salant, S., and Fitzgibbons, P. J. (2004). ‘‘Effects of stimulus and noise rate variability on speech perception by younger and older adults,’’ *J. Acoust. Soc. Am.* **115**(4), 1808–1817.

- Grieser, D. A. L., and Kuhl, P. K. (1988). "Maternal speech to infants in a tonal language: Support for universal prosodic features in motherese," *Dev. Psychol.* **24**(1), 14–20.
- Hualde, J., and Chitoran, I. (2003). "Explaining the distribution of hiatus in Spanish and Romanian," in *Proc. Int. Conf. Phonetic Sciences*, Barcelona, pp. 1683–1686.
- Jokinen, E., Remes, U., and Alku, P. (2016). "The use of read versus conversational Lombard speech in spectral tilt modeling for intelligibility enhancement in near-end noise conditions," in *Proc. Interspeech*, pp. 2771–2775.
- Junqua, J.-C. (1993). "The Lombard reflex and its role on human listeners and automatic speech recognizers," *J. Acoust. Soc. Am.* **93**(1), 510–524.
- Koutsogiannaki, M., and Stylianou, Y. (2016). "Modulation enhancement of temporal envelopes for increasing speech intelligibility in noise," in *Interspeech 2016*, pp. 2508–2512.
- Lu, Y., and Cooke, M. (2008). "Speech production modifications produced by competing talkers, babble, and stationary noise," *J. Acoust. Soc. Am.* **124**(5), 3261–3275.
- Mendoza, E., Carballo, G., Cruz, A., Fresneda, M. D., Muoz, J., and Marrero, V. (2003). "Temporal variability in speech segments of Spanish: Context and speaker related differences," *Speech Commun.* **40**, 431–447.
- Miller, G. A., and Licklider, J. C. (1950). "The intelligibility of interrupted speech," *J. Acoust. Soc. Am.* **22**, 167–173.
- Moore, B. C. J., Glasberg, B. R., Plack, C. J., and Biswas, A. K. (1988). "The shape of the ear's temporal window," *J. Acoust. Soc. Am.* **83**(7–8), 1102–1116.
- Moreno, A., Poch, D., Bonafonte, A., Lleida, E., Llisterri, J., Marino, J. B., and Nadeu, C. (1993). "Albayzín speech database: Design of the phonetic corpus," in *Eurospeech*, Berlin, Germany, pp. 175–78.
- Naylor, G. (2016). "Theoretical issues of validity in the measurement of aided speech reception threshold in noise for comparing nonlinear hearing aid systems," *J. Am. Acad. Audiol.* **27**, 504–514.
- Nejime, Y., and Moore, B. C. J. (1998). "Evaluation of the effect of speech-rate slowing on speech intelligibility in noise using a simulation of cochlear hearing loss," *J. Acoust. Soc. Am.* **103**(1), 572–576.
- Oxenham, A. J., and Simonson, A. M. (2009). "Masking release for low- and high-pass-filtered speech in the presence of noise and single-talker interference," *J. Acoust. Soc. Am.* **125**, 457–468.
- Patterson, R. D., Holdsworth, J., Nimmo-Smith, I., and Rice, P. (1988). "SVOS Final Report: The Auditory Filterbank, Technical Report 2341 (MRC Applied Psychology Unit, Cambridge, UK).
- Picheny, M. A., Durlach, N. I., and Braida, L. D. (1985). "Speaking clearly for the hard of hearing. I: Intelligibility differences between clear and conversational speech," *J. Speech Hear. Res.* **28**, 96–103.
- Pisoni, D. B., Bernacki, R. H., Nusbaum, H. C., and Yuchtman, M. (1985). "Some acoustic-phonetic correlates of speech produced in noise," in *ICASSP*, Tampa, FL, pp. 1581–1584.
- Pittman, A. L., and Wiley, T. L. (2001). "Recognition of speech produced in noise," *J. Speech Lang. Hear. Res.* **44**(3), 487–496.
- Rothauser, E. H., Chapman, W. D., Guttman, N., Hecker, M. H. L., Nordby, K. S., Silbiger, H. R., Urbanek, G. E., Weistock, M., McGee, V. E., Pacht, U. P., and Voiers, W. D. (1969). "IEEE recommended practice for speech quality measurements," *IEEE Trans. Audio Acoust.* **17**, 225–246.
- Sankowska, J., García Lecumberri, M. L., and Cooke, M. (2011). "Interaction of intrinsic vowel and consonant durational correlates with foreigner directed speech," *Poznań Studies Contemp. Linguist.* **47**, 109–119.
- Sauert, B., and Vary, P. (2006). "Near end listening enhancement: Speech intelligibility improvement in noisy environments," in *Proc. ICASSP*, Toulouse, France, pp. 493–496.
- Schepker, H., Rennie, J., and Doclo, S. (2013). "Improving speech intelligibility in noise by sii-dependent preprocessing using frequency-dependent amplification and dynamic range compression," in *Proc. Interspeech*, pp. 3577–3581.
- Skowronski, M. D., and Harris, J. G. (2006). "Applied principles of clear and Lombard speech for automated intelligibility enhancement in noisy environments," *Speech Commun.* **48**(5), 549–558.
- Song, J. Y., Demuth, K., and Morgan, J. (2010). "Effects of the acoustic properties of infant-directed speech on infant word recognition," *J. Acoust. Soc. Am.* **128**(1), 389–400.
- Stilp, C., and Kluender, K. (2010). "Cochlea-scaled entropy, not consonants, vowels, or time, best predicts speech intelligibility," *Proc. Natl. Acad. Sci. U.S.A.* **107**(27), 12387–12392.
- Summers, W. V., Pisoni, D. B., Bernacki, R. H., Pedlow, R. I., and Stokes, M. A. (1988). "Effects of noise on speech production: Acoustic and perceptual analyses," *J. Acoust. Soc. Am.* **84**(3), 917–928.
- Taal, C. H., Jensen, J., and Leijon, A. (2013). "On optimal linear filtering of speech for near-end listening enhancement," *IEEE Signal Proc. Lett.* **20**(3), 225–228.
- Tang, Y., and Cooke, M. (2012). "Optimised spectral weightings for noise-dependent speech intelligibility enhancement," in *Proc. Interspeech*, Portland, OR, pp. 955–958.
- Uchanski, R. M. (2005). "Clear speech," in *The Handbook of Speech Perception*, edited by D. B. Pisoni, and R. E. Remez (Blackwell, Oxford, UK), pp. 207–235.
- Valentini-Botinhao, C., Yamagishi, J., and King, S. (2012). "Mel cepstral coefficient modification based on the Glimpse Proportion measure for improving the intelligibility of HMM-generated synthetic speech in noise," in *Proc. Interspeech*, Portland, OR, pp. 631–634.
- Yoo, S. D., Boston, J. R., El-Jaroudi, A., Li, C.-C., Durrant, J. D., Kovacyk, K., and Shaiman, S. (2007). "Speech signal modification to increase intelligibility in noisy environments," *J. Acoust. Soc. Am.* **122**(2), 1138–1149.
- Zorila, T., Kandia, V., and Stylianou, Y. (2012). "Speech-in-noise intelligibility improvement based on spectral shaping and dynamic range compression," in *Proc. Interspeech*, pp. 635–638.
- Zorila, T.-C., and Stylianou, Y. (2015). "A fast algorithm for improved intelligibility of speech-in-noise based on frequency and time domain energy reallocation," in *Proc. Interspeech*, pp. 60–64.