

Spectral completion of partially masked sounds

Josh H. McDermott* and Andrew J. Oxenham

Department of Psychology, University of Minnesota, N640 Elliott Hall, 75 East River Road, Minneapolis, MN 55455-0344

Edited by Eric I. Knudsen, Stanford University School of Medicine, Stanford, CA, and approved February 15, 2008 (received for review November 29, 2007)

Natural environments typically contain multiple sound sources. The sounds from these sources frequently overlap in time and often mask each other. Masking could potentially distort the representation of a sound's spectrum, altering its timbre and impairing object recognition. Here, we report that the auditory system partially corrects for the effects of masking in such situations, by using the audible, unmasked portions of an object's spectrum to fill in the inaudible portions. This spectral completion mechanism may help to achieve perceptual constancy and thus aid object recognition in complex auditory scenes.

auditory objects | auditory scene analysis | perceptual organization | segmentation | segregation

The presence of multiple sound sources is a routine occurrence in the natural world but poses a challenge to the auditory system, which must separate each source from the sum of the source waveforms (1–3). This challenge is compounded by the frequent occurrence of masking (4), in which sounds of interest are partially obscured by other sufficiently loud sounds. Masking introduces distortions that could impair the identification of a sound and generally alter how it is heard. Auditory scene analysis is thus believed to entail compensatory mechanisms to help infer the true characteristics of a sound, i.e., those that would be heard in the absence of masking.

Thus far, the primary documented means for this has been the so-called “continuity illusion.” It has long been known that sounds interrupted by brief masking noises are heard to continue through them despite the physical disruption caused by the masker (5–8). The effect occurs for stimuli ranging from tones to speech syllables (9); the masking noise bursts used in laboratory conditions mimic the effect of handclaps, coughs, and other common brief masking sounds. Although the mechanisms of this effect remain poorly understood (10–12), it presumably functions to produce perceptual continuity in conditions where the original source is likely to have been continuous, even though the stimulus entering the ear is not.

Many environments present a different challenge, because of sounds that are extended in time, such as those produced by an office fan, a river, or chatter in a crowded room. Because such background sounds are temporally extended, there may be little disruption of the temporal continuity of sounds of interest. However, masking can nonetheless occur, and because masking sounds are often not spectrally uniform, they have the potential to obscure some portions of an object's spectrum but not others. If uncorrected, such masking could lead to perceptual distortions. In this article, we explore whether the auditory system might correct for these distortions by using audible portions of an object's spectrum to infer the portions that might likely be masked.

We studied the simple case in which two sounds overlap in time and frequency and therefore have the potential to mask each other. Consider the stimulus of Fig. 1*a*, depicted with a schematic spectrogram. Energy is present in low-, middle-, and high-frequency bands, but the high and low bands start later and end earlier than the middle band. The different onset and offset times would be expected to produce the perception of two distinct sounds, and indeed this is what listeners report hearing: a long narrowband noise overlapped by a second, briefer noise

burst. The stimulus renders the precise characteristics of the second sound ambiguous. It could simply consist of the high and low bands alone, because these could be segmented from the middle band by virtue of their delayed onset. However, the stimulus leaves open a second possibility—that the briefer sound contains energy in the middle band that is masked by the longer sound. The continuous nature of many natural sound spectra might favor such an interpretation, but it remains to be seen whether listeners actually hear sounds in this way.

Results

Experiment 1. We used a matching task in which subjects heard a standard stimulus (e.g., Fig. 1*a*) and then adjusted the middle band of a subsequent comparison stimulus (Fig. 1*b*). The standard typically was designed to yield the percept of two sounds described above, one long and one short. Subjects were instructed to direct their attention to the shorter sound, termed the “target.” The comparison stimulus was designed to be heard as a single sound of the same duration as the target. Subjects were instructed to make the comparison sound as similar as possible to the target. The high and low bands of the comparison stimulus were fixed to be identical to those in the standard, and subjects adjusted the level of the middle band to create a perceptual match. If the auditory system infers the target sound to contain energy in the middle band, subjects' matches ought to reflect this. For clarity, we will refer to the long middle band of noise as the “masker,” and the high- and low-frequency noise bursts as the “tabs.”

To first confirm that subjects could accurately perform the task, we measured their matches in two control conditions in which the masker was absent (Fig. 1*c*, i and ii). As expected, when presented with just the tabs, subjects set the comparison middle band to very low levels, in the neighborhood of the detection threshold for such stimuli (13). In a second condition, the middle band of the standard was high in level (spectrum level of 30 dB re: 20 μ Pa) but was the same duration as the tabs, such that a single brief sound was perceived. Subjects' matches were again close to veridical (compare with filled circle in Fig. 1*c*), suggesting that they were able to do the task with reasonable accuracy.

When the masker was combined with the tabs in the condition of interest (Fig. 1*c*, iii), subjects adjusted the comparison middle band far above detection thresholds, indicating that the target seemed to contain middle band frequencies—the tabs by themselves were an inadequate match to what subjects heard. This effect depended critically on the presence of masker energy at the appropriate location. When the masker contained a large spectral gap or had a temporal gap coincident with the onset and offset of the high and low bands (Fig. 1*c*, iv and v), subjects assigned much less energy to the middle band. For additional results, see [supporting information \(SI\) Results](#) and [Fig. S1](#).

Author contributions: J.H.M. and A.J.O. designed research; J.H.M. performed research; J.H.M. analyzed data; and J.H.M. and A.J.O. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

*To whom correspondence should be addressed. E-mail: joshmcd@umn.edu.

This article contains supporting information online at www.pnas.org/cgi/content/full/0711291105/DCSupplemental.

© 2008 by The National Academy of Sciences of the USA

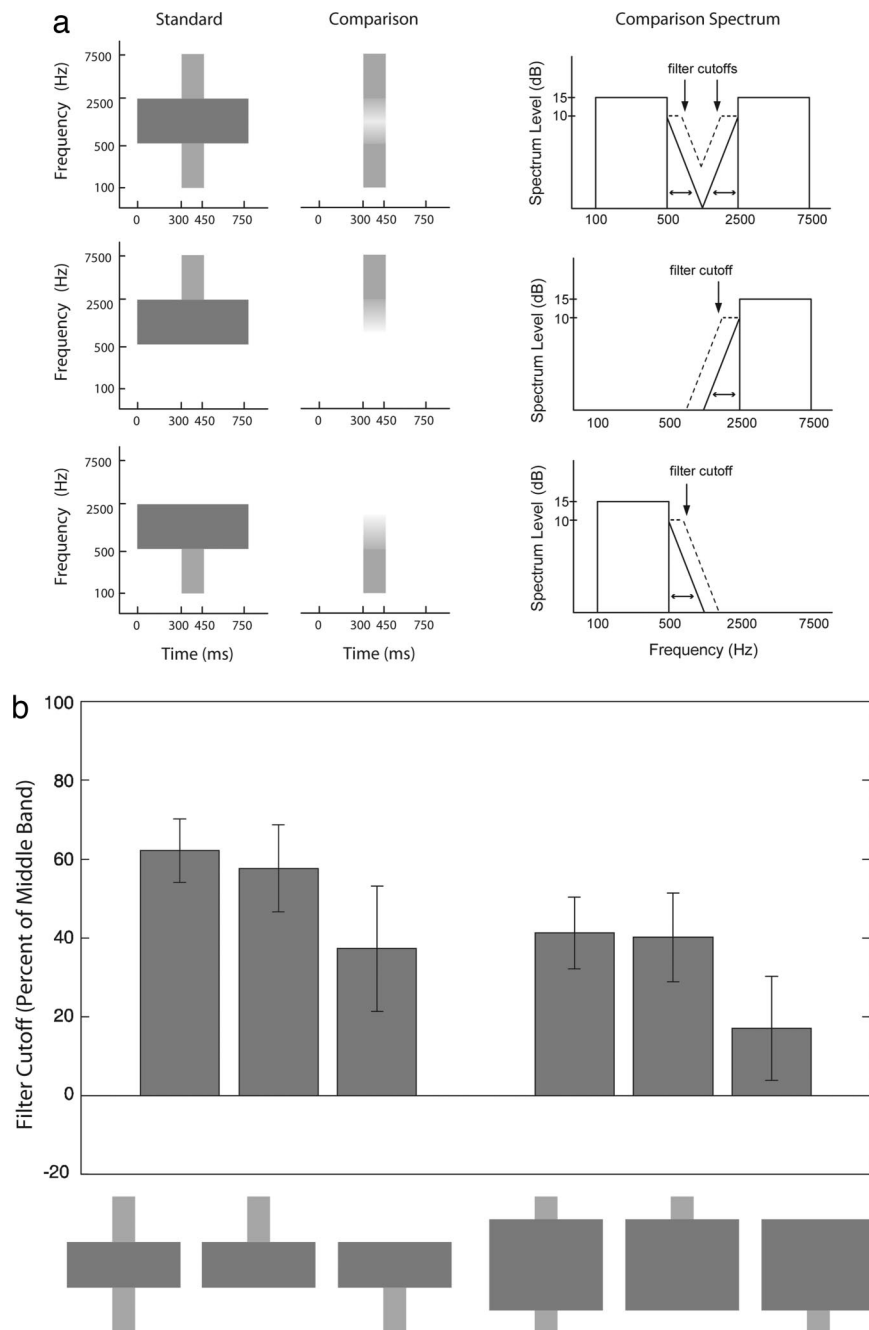


Fig. 4. Effect of high and low inducing elements alone. (a) Stimuli; only the stimuli with narrow masker bandwidths are shown. The adjustable middle band noise was generated with a filter whose cutoff was adjusted by subjects. Two settings are shown in the schematic spectra on the right (solid, low percentage; dashed, higher percentage, with arrows denoting cutoffs for higher percentage setting). (b) Mean filter cutoff settings (six subjects).

used a second-order filter with a shallow roll-off to generate the matching noise, so even when the cutoff is at zero (i.e., when it is at the borders of the middle band), a substantial amount of noise is added to the middle band.

We again observed a main effect of masker bandwidth [Fig. 4b, left vs. right; $F(1, 5) = 13.3$, $P = 0.015$], but found no significant effect of the tab configuration [$F(2, 10) = 1.51$, $P = 0.266$], and no interaction [$F(2, 10) = 0.04$, $P = 0.97$]. There is a nonsignificant trend for more completion to occur for high-frequency tabs than for low, but it is clear that the effect persists with a single tab alone. The effect of both high and low tabs at once is not appreciably more than the sum of the effects of the

high and low tabs, because the cutoff settings are similar in all three conditions.

Experiment 5: Completion of Complex Tones. Similar effects were also observed with harmonic sounds more similar to those found in speech and music. When a subset of the components of a harmonic complex tone was presented above or below bandpass noise (Fig. 5a), most subjects reported the perceived brightness (16) to be altered. Masking of the components by the noise predicts that the masker should raise the brightness of the high tone and lower that of the low tone, because the audibility of components close to the masker would be reduced. In fact, we observed the masking noise to have the opposite effect, consis-

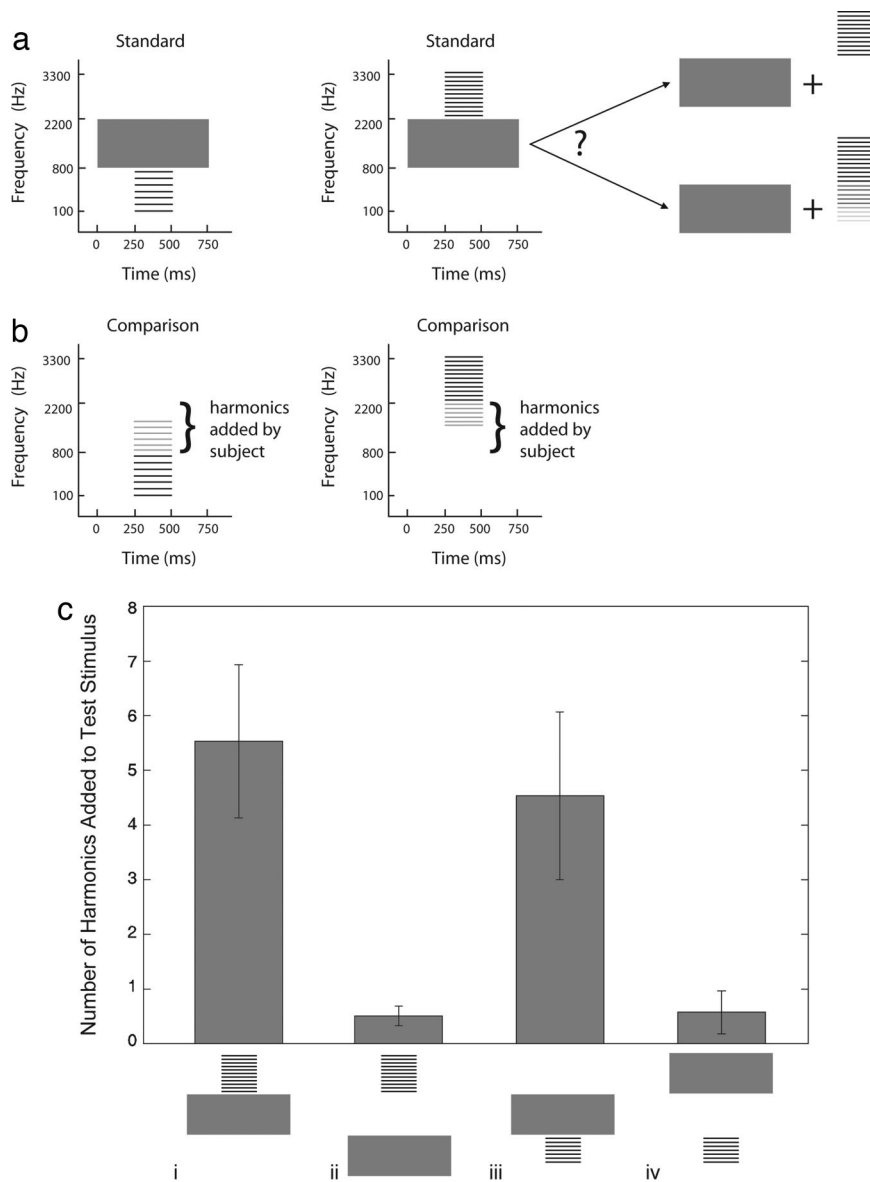


Fig. 5. Spectral completion of complex tones. (a) Schematics of standard with inducing harmonics in upper spectral region. Two possible perceptual interpretations are shown to the right. Spectral completion might cause subjects to hear harmonics in the spectral region of the masker, in this case reducing the tone brightness by lowering the perceived spectral centroid. (b) Schematic comparison stimulus for the conditions with harmonics in the upper frequency band. Subjects adjusted the number of harmonics added to those present in the standard. (c) Mean number of harmonics added to masked tones (six subjects). Stimuli are not drawn exactly to scale in frequency domain.

tent with the possibility that the auditory system infers frequency components that would be obscured by the masker. To quantify this, we had subjects perform a matching task in which they added low-amplitude harmonics to a comparison stimulus to make it resemble the sound of the tone in the standard (Fig. 5*b*). When masking noise was presented in the spectral region adjacent to that of the tone burst, subjects added harmonics to the middle band of the comparison stimulus (the region of the masker of i and iii), whereas noise bursts of equal amplitude presented in a nonadjacent spectral region had no such effect [Fig. 5*c*; $F(1, 5) = 188.259$, $P < 0.0001$; data square-root transformed to normalize variance].

Discussion

These experiments suggest that the auditory system uses unmasked spectral regions of sounds to infer the regions that are

likely to have been masked. The effect occurs for both tonal and nontonal sounds and seems to respect the physical constraints of masking, positing only as much spectral energy as is consistent with the masker level. The resulting perceived spectral content of sounds presented adjacent to potential maskers is opposite to that predicted by conventional masking.

The mechanism documented here seems to complement the classic “continuity” effect (5). Whereas the continuity mechanism links segments of sound across time to compensate for temporal disruption by intermittent masking sounds, the proposed spectral completion process compensates for masking of part of the spectrum by a continuous masker, via completion in the frequency domain. Spectral completion would seem to function primarily to achieve a faithful representation of an object’s spectrum during masking, the main result of which would be to promote timbre constancy. In contrast, continuity in

time does not alter timbre but does affect the perception of temporal structure, which our proposed process leaves unaffected. The two effects thus appear complementary, helping to solve different problems for the auditory system.

These results have interesting implications for theories of auditory scene analysis. Standard scene analysis models posit that onset cues are used to assign spectral energy to the various sound sources in a scene (1–3). The effects described here suggest that under conditions in which masking is likely to occur, the auditory system assigns spectral energy to sound sources even in the absence of onset cues in the assigned frequency channels, by extrapolating from adjacent spectral regions that themselves contain onsets. Previous studies have shown that adding noise to spectral gaps in speech sounds can enhance intelligibility (9, 18, 19); our results suggest that this may reflect a spectral completion process. Such a process cannot, of course, fully circumvent the effects of masking, but it may help to reduce the distortions in perception that would otherwise occur from partial masking of the spectrum.

Materials and Methods

General. A single trial within an iterative run consisted of a presentation of the standard and comparison stimuli for a given condition. The standard was fixed within a run; after each iteration, subjects had the option of adjusting the level of the middle band in the comparison stimulus. The starting level of the middle band was chosen randomly between -10 and 30 dB (spectrum level re: 20 μ Pa). Iterations were self-paced and continued until a subject determined that they had achieved a satisfactory match, at which point they clicked a button to move to the next run. The level on the last iteration of each run was stored as the matching level for that run. The order of presentation of the conditions in an experiment was randomized.

All subjects had normal hearing, as defined by pure-tone thresholds of 20 dB hearing loss or less at octave frequencies between 250 and 8000 Hz, and did not report any history of hearing disorders. Subjects (18 – 30 years of age) began by completing a session's worth of practice runs (typically 10 runs per condition) that were not included in the data analysis. Some subjects declined to return for the experimental sessions or did not complete the full allotment of runs (20 runs per condition in all experiments) and were not included in the analysis.

Stimuli were generated by combining band-limited Gaussian noise bursts. Each burst was generated in the spectral domain within a single buffer, setting all magnitude coefficients outside the spectral pass band to zero and performing an inverse fast Fourier transform. The pass bands of lower tab, masker, and upper tab extended from 100 to 500 Hz, 500 to 2500 Hz, and 2500 to 7500 Hz, respectively (Experiments 1–4). The upper-tab bandwidth was narrower than that of the lower tab on a log scale to more closely approach equal loudness of the tabs. The total masker duration was 750 ms, and the duration of the upper and lower tabs was 150 ms. In the standard stimuli, the tabs started 300 ms after the onset of the masker. All stimuli were gated on and off with 10 -ms raised-cosine (Hanning) ramps. The time interval between the end of the standard and beginning of the comparison in each iteration of a trial was 300 ms.

Sounds were generated digitally and played out by a LynxStudio Lynx22

24-bit D/A converter at a sampling rate of 48 kHz. The sounds were then presented diotically to subjects through Sennheiser HD580 headphones.

Experiment 1 (Fig. 1). The spectrum level of the tabs and masker in their pass bands was 20 dB (re: 20 μ Pa). The spectrum level of the middle band of the standard in Fig. 1c, part ii, was 30 dB. The spectral gap in the stimulus of Fig. 1c, part iv, extended from 600 to $2,080$ Hz (chosen such that the long masker bands were equally wide on a log scale); the spectrum level of the masker bands was raised so that the overall level of the masker was equated to that of the other conditions.

Experiment 2 (Fig. 2). The spectrum levels of the tabs and masker were varied in opposite directions across conditions (5 and 35 , 10 and 30 , 15 and 25 , 20 and 20 , 25 and 15 , 30 and 10 dB, tabs and masker, respectively).

Experiment 3 (Fig. 3). Part a: The upper border of the lower tab and the lower border of the upper tab were altered so as to introduce gaps or vary the masker bandwidth. Altered borders of tabs: 170 and $5,200$ Hz in i and v, 290 and $3,600$ Hz in ii and iv. Spectrum level of tabs and masker in their pass bands was 20 dB. Part b: Both borders of both tabs were shifted so as to maintain constant bandwidth on an ERB scale, of 3 ERBs [lower tab: 100 and 226 Hz (i and v), 226 and 400 Hz (ii and iv), and 400 and 640 Hz (iii); upper tab: $5,724$ and $8,000$ Hz (i and v), $4,077$ and $5,724$ Hz (ii and iv), and $2,886$ and $4,077$ Hz (iii)]. The masker borders were 640 and $2,886$ Hz in i, ii, and iii and otherwise were equal to the upper cutoff of the lower tab and the lower cutoff of the upper tab.

Experiment 4 (Fig. 4). The spectrum level of the masker and tabs was 25 and 15 dB, respectively. The adjustable middle band noise in the comparison stimulus was generated with a second-order Butterworth filter, the cutoff frequency of which was adjusted by subjects. Filter cutoffs were defined as the point of 3 -dB attenuation; because of the shallow roll-off, subjects were allowed to adjust cutoffs to values outside the middle band (in which case, the percentage of the band filled was negative). The spectrum level of this noise where it was unattenuated by the filter was 10 dB. The band borders of the stimuli with narrower maskers (on the left of Fig. 4c) were as in Experiment 1; for the stimuli with broader maskers, they were 100 , 290 , $5,200$, and $7,500$ Hz.

Experiment 5 (Fig. 5). The high and low tones were composed of evenly spaced harmonics from $2,300$ to $3,300$ Hz and 100 to 700 Hz, respectively, in steps of 100 Hz. The high tones had more harmonics and were at a higher level (50 dB vs. 40 dB SPL per harmonic for the low tones) to make them approximately as loud as the low tones. The harmonics added to the middle band (the band occupied by the noise masker in the standard stimulus of i and iii) extended up or down from the highest/lowest harmonic of the low and high tone bursts, respectively, with the same spacing. The initial number of harmonics in the middle band was chosen randomly between 1 and 11 . They were 10 dB lower in level than the inducing harmonics. Tone bursts were 250 ms in length; noise maskers were 750 ms. The maskers extended from 800 to $2,200$ Hz (i and iii), from 100 to 800 Hz (ii), and from $2,200$ to $5,000$ Hz (iv). The masker spectrum level was 35 dB in i and iii; in ii and iv, the maskers were scaled such that the overall level was the same across conditions.

ACKNOWLEDGMENTS. We thank Christophe Micheyl, Tali Sharot, and Jonathan Winawer for helpful comments on the manuscript. This work was supported by National Institutes of Health Grant R01 DC 07657.

- Bregman AS (1990) *Auditory Scene Analysis: The Perceptual Organization of Sound* (MIT Press, Cambridge, MA).
- Carlyon RP (2004) How the brain separates sounds. *Trends Cognit Sci* 8:465–471.
- Darwin CJ (1997) Auditory grouping. *Trends Cognit Sci* 1:327–333.
- Moore BCJ (1995) Frequency analysis and masking. *Handbook of Perception and Cognition*, ed Moore BCJ (Academic, Orlando, FL), Vol 6, pp 161–205.
- Warren RM (1970) Perceptual restoration of missing speech sounds. *Science* 167:392–393.
- Houtgast T (1972) Psychophysical evidence for lateral inhibition in hearing. *J Acoust Soc Am* 51:1885–1894.
- Dannenbring GL (1976) Perceived auditory continuity of alternately rising and falling FM sweeps. *Can J Psychol* 30: 99–114.
- Carlyon RP, et al. (2004) Auditory processing of real and illusory changes in frequency modulation (FM) phase. *J Acoust Soc Am* 116:3629–3639.
- Warren RM (1999) *Auditory Perception: A New Analysis and Synthesis* (Cambridge Univ Press, Cambridge, UK).
- Micheyl C, et al. (2003) The neurophysiological basis of the auditory continuity illusion: A mismatch negativity study. *J Cognit Neurosci* 15:747–758.
- Riecke L, et al. (2007) Hearing illusory sounds in noise: Sensory-perceptual transformations in primary auditory cortex. *J Neurosci* 27:12684–12689.
- Petkov CI, O'Connor KN, Sutter ML (2007) Encoding of illusory continuity in primary auditory cortex. *Neuron* 54:153–165.
- Glasberg BR, Moore BCJ (1990) Derivation of auditory filter shapes from notched-noise data. *Hear Res* 47:103–138.
- Schacknow PN, Raab DH (1976) Noise-intensity discrimination: effects of bandwidth conditions and mode of masker presentation. *J Acoust Soc Am* 60:893–905.
- Shipley TF, Kellman PJ (1992) Strength of visual interpolation depends on the ratio of physically specified to total edge length. *Percept Psychophys* 52:97–106.
- Stevens JC, Hall JW (1966) Brightness and loudness as a function of stimulus duration. *Percept Psychophys* 1: 319–327.
- von Bismark G (1974) Sharpness as an attribute of the timbre of steady sounds. *Acustica* 30:159–172.
- Shriberg EE, Perceptual restoration of filtered vowels with added noise. *Language Speech* 35:127–136.
- Warren RM, et al. (1997) Spectral restoration of speech: intelligibility is increased by inserting noise in spectral gaps. *Percept Psychophys* 59:275–283.