

Central Lancashire Online Knowledge (CLoK)

Title	When softer sounds are more distracting: Task-irrelevant whispered speech causes disruption of serial recall
Type	Article
URL	https://clock.uclan.ac.uk/53619/
DOI	
Date	2024
Citation	Kattner, Florian, Focker, Julia, Moshona, Cleopatra Christina and Marsh, John Everett (2024) When softer sounds are more distracting: Task-irrelevant whispered speech causes disruption of serial recall. Journal of the Acoustical Society of America (JASA). ISSN 0001-4966
Creators	Kattner, Florian, Focker, Julia, Moshona, Cleopatra Christina and Marsh, John Everett

It is advisable to refer to the publisher's version if you intend to cite from the work.

For information about Research at UCLan please go to <http://www.uclan.ac.uk/research/>

All outputs in CLoK are protected by Intellectual Property Rights law, including Copyright law. Copyright, IPR and Moral Rights for the works on this site are retained by the individual authors and/or other copyright owners. Terms and conditions for use of this material are defined in the <http://clock.uclan.ac.uk/policies/>

When softer sounds are more distracting: Task-irrelevant whispered speech causes disruption of serial recall

Florian Kattner,¹ Julia Föcker,² Cleopatra Christina Moshona,³ and John E. Marsh⁴

¹*Institute for Mind, Brain and Behavior, Health and Medical University,
Schiffbauergasse 14, 14467 Potsdam, Germany^a*

²*College of Health and Science, School of Psychology, University of Lincoln,
Brayford Pool, Lincoln, LN6 7TS, United Kingdom*

³*Engineering Acoustics, Institute of Fluid Dynamics and Technical Acoustics,
Technische Universität Berlin, Einsteinufer 25, 10587 Berlin,
Germany*

⁴*School of Psychology and Humanities, University of Central Lancashire, Preston,
PR1 2HE, United Kingdom*

(Dated: 2 November 2024)

Abstract

Two competing accounts propose that the disruption of short-term memory by irrelevant speech arises either due to interference-by-process (e.g., changing-state effect) or attentional capture, but it is unclear how whispering affects the irrelevant speech effect. According to the interference-by-process account, whispered speech should be less disruptive due to its reduced periodic spectro-temporal fine structure and lower amplitude modulations. In contrast, the attentional account predicts more disruption by whispered speech, possibly via enhanced listening effort in the case of a comprehended language. In two experiments, voiced and whispered speech (spoken sentences or monosyllabic words) were presented while participants memorized the order of visually presented letters. In both experiments, a changing-state effect was observed regardless of the phonation (sentences produced more disruption than 'steady-state' words). Moreover, whispered speech (lower fluctuation strength) was more disruptive than voiced speech when participants understood the language (Exp. 1), but not when the language was incomprehensible (Exp. 2). The results suggest two functionally distinct mechanisms of auditory distraction: While changing-state speech causes automatic interference with seriation processes regardless of its meaning or intelligibility, whispering appears to contain cues that divert attention from the focal task primarily when presented in a comprehended language, possibly via enhanced listening effort.

^aCorrespondence to: florian.kattner@hmu-potsdam.de

21 **I. INTRODUCTION**

22 Most readers will have experienced disruption to their cognitive performance in the pres-
23 ence of task-irrelevant background sound even when the focal task information is in a differ-
24 ent modality (e.g., visual) and therefore cannot be attributed to interference at the sensory
25 level (e.g., perceptual masking). Instead it must emerge from an interaction between visual
26 and auditory processing at a level beyond the sensory organs. One well-studied example
27 of auditory distraction is the disruption of verbal-serial short-term memory produced by
28 task-irrelevant speech (Colle and Welsh, 1976; Salamé and Baddeley, 1982). In this irrele-
29 vant sound paradigm, participants are asked to recall a series of usually visually-presented
30 digits or words while being presented with different types of sound via headphones that they
31 are instructed to deliberately ignore. Immediate or delayed visual-verbal serial recall accu-
32 racy is usually lower when task-irrelevant speech is presented during encoding or retention
33 of the items, compared to silence, continuous noise, or instrumental background music (in
34 particular when the notes are played “legato”, i.e., smoothly connected without gaps of si-
35 lence; Ellermeier and Zimmer, 1997; Salamé and Baddeley, 1989; Schlittmeier *et al.*, 2008),
36 regardless of the volume of irrelevant speech (Ellermeier and Hellbrück, 1998). However,
37 distraction of visual-verbal serial recall is not restricted to speech or “speech-like” mate-
38 rial (e.g., music; Salamé and Baddeley, 1989), as stimuli sufficiently unlike speech such as
39 spectro-temporally varying tones (Jones and Macken, 1993) or pitch glides randomly inter-
40 rupted with quiet (Jones *et al.*, 1993) have also been found to interfere with the serial order
41 retention of to-be-remembered items. However, the magnitude of the disruptive effect of

42 non-speech sound (e.g., music, varying tones, or interrupted pitch glides) was found to be
43 significantly lower than that of irrelevant speech (see effect sizes in [Ellermeier and Zimmer,](#)
44 [2014](#), Table 1).

45 The disruptive impact of irrelevant speech is observed even if it is presented only during
46 a retention period after encoding of the visually-presented items. Therefore, the disruption
47 is not due to interference in the encoding of digits, but occurs at a later stage of processing
48 within memory ([Miles *et al.*, 1991](#)). Within the context of short-term memory, two broad
49 mechanisms have been proposed to account for auditory distraction. According to the
50 ‘interference-by-process account’ (which can be considered a generalization of the ‘object-
51 oriented episodic record account’, [Jones *et al.*, 1996](#); [Marsh *et al.*, 2009](#)), processing of
52 task-irrelevant sound produces interference with cognitive processes that are demanded by
53 the focal task. One prototypical example of such interference is the changing-state effect,
54 which refers to the observation that spectro-temporally varying sound (e.g., free-running
55 speech or random sequences of syllables or tones) gives rise to the automatic formation of
56 an ordered auditory sequence (as part of auditory scene analysis, [Bregman, 1990](#)) which
57 then interferes with deliberate serial-order processing of to-be-remembered information. In
58 line with this assumption, it has been found that changing-state sequences consisting of
59 spectro-temporally varying acoustical tokens (thus conveying irrelevant order information)
60 are more disruptive than steady-state repetitions of a single acoustical item, regardless of
61 whether the sequences comprise speech or non-speech materials (e.g., [Hadlington *et al.*, 2004](#);
62 [Jones and Macken, 1993](#); [Jones *et al.*, 1993](#); [Tremblay *et al.*, 2000](#), but see [LeCompte *et al.*,](#)
63 [1997](#)). Moreover, this well-established ‘changing-state effect’ ([Jones *et al.*, 1992](#)) seems to

64 interfere primarily in tasks that require serial-order processing or with the performance of
65 participants who report using serial rehearsal for item retention, whereas often no changing-
66 state effect is found with non-serial memory tasks such as the missing item task or in
67 a mental arithmetic task (Beaman and Jones, 1997; Campbell *et al.*, 2002; Hughes and
68 Marsh, 2020; Jones and Macken, 1993; Kattner *et al.*, 2023). Importantly, according to the
69 interference-by-process account, auditory distraction in a serial recall task should depend
70 primarily on the acoustical profile of the irrelevant sound (e.g., the proportion of changes
71 in rhythm, frequency, or amplitude), and it has been discussed whether psychoacoustical
72 metrics such as the degree of amplitude or frequency modulation, ‘fluctuation strength’
73 , or spectral detail may be useful predictors of the changing-state effect (Ellermeier and
74 Zimmer, 2014; Schlittmeier *et al.*, 2012). It is worth noting that the degree of distraction
75 imposed by the changing-state sound may also depend on certain speech-specific properties
76 of the sound. For instance, it has been found that artificial sinewave speech containing three
77 formants can be as disruptive as natural speech, but a temporal reversal of the first two
78 formants (i.e., degrading the formant transitions that are required to identify ‘consonants’
79 in sine wave speech) reduced the degree of distraction considerably (Viswanathan *et al.*,
80 2014). Similarly, Dorsi *et al.* 2018 found that if the spectral detail of the irrelevant speech is
81 reduced by decreasing the number of vocoder bands, the distraction caused by task-irrelevant
82 speech diminishes (see also Ellermeier *et al.*, 2015). Taken together, these findings suggest
83 that speech-specific properties and signal fidelity (i.e. the internal properties of a signal,
84 including its structural details and relationships) may play a cardinal role in modulating the
85 effects of task-irrelevant speech.

86 In contrast, an alternative ‘unitary attentional account’ supposes that irrelevant sounds
87 divert attentional or cognitive resources from the focal task (Bell *et al.*, 2008, 2012; Cowan,
88 1995). Specifically, certain types of irrelevant sound (e.g., speech, acoustical changes, un-
89 expected events, or otherwise meaningful sounds) are assumed to capture attention and
90 produce unspecific disruption to any attention-demanding task. By this approach, a sound
91 may capture attention either because it cannot be predicted based on previous stimulation
92 (random acoustical changes or an auditory oddball in a regular sequence, Eimer *et al.*, 1996),
93 or because semantic or syntactic properties of the sound indicate enhanced relevance to the
94 individual (e.g., one’s own name or an emotional word Röer *et al.*, 2013, 2017a). According
95 to this account, the degree of disruption should not depend on the exact cognitive processes
96 demanded by the focal task (e.g., retention of order), but it should vary as a function of the
97 perceptual load and/or the working memory capacity available to the participant. Indeed,
98 it has been reported that the disruptive effect of a deviant (unexpected) sound in a regu-
99 lar sequence is more pronounced in, or restricted to, individuals with low working-memory
100 capacity and conditions of low task-encoding load (Hughes *et al.*, 2013; Hughes and Marsh,
101 2019; Marsh *et al.*, 2018; Sörqvist, 2010)(but see Körner *et al.*, 2017; Labonté *et al.*, 2022).
102 However, in contrast to a unitary attentional account, the disruptive effect of other types of
103 sounds, in particular changing-state sound, does not seem to depend on task load and the
104 individuals’ working memory capacity (Hughes *et al.*, 2013) – though it might be sensitive
105 to the listeners’ auditory processing and/or selective attention (cf. reduced distraction in
106 blind individuals Kattner and Ellermeier, 2014; Kattner *et al.*, 2024).

107 To account for such findings, a duplex-mechanism account of auditory distraction has
108 been proposed, assuming that interference-by-process and attentional capture may be two
109 functionally distinct mechanisms that can produce task disruption (Hughes, 2014; Hughes
110 *et al.*, 2005b, 2007). That is, irrelevant sound may either produce interference with specific
111 cognitive processes that are demanded by the focal task (Kattner, 2024; Marsh *et al.*, 2009,
112 e.g., changing-state sound interferes with a seriation process, and semantic properties of
113 irrelevant sound may interfere with semantic organization; cf.) or it may capture attention
114 due to its unpredictability or meaningfulness and cause unspecific disruption (assuming that
115 sufficient attentional/cognitive resources are available to process the sound).

116 Whispered speech is an interesting stimulus to test the functional dissociation between
117 interference-by-process (i.e., interference with seriation) and attentional capture. In contrast
118 to voiced (modal) phonation, vocal cord vibration, periodic glottal excitation and harmonic
119 structure are completely absent in whispered speech, due to its distinct production mech-
120 anism. The glottis is abducted, except for a small triangular opening in the cartilaginous
121 portion (Laver, 1994). The pulmonic airstream forced through this narrow gap has a hiss-
122 ing, noise-like quality, produced by turbulence from the friction of the air around the larynx
123 (Eckert and Laver, 1994). Consequently, whispered speech is dominated by strong aperiodic
124 energy. It is further characterized by a notable decrease in vowel amplitude, typically by
125 about 20–25 dB, flatter spectral slopes and an upwards shift of formant frequencies, affecting
126 vowel quality and intelligibility (Ito *et al.*, 2005). These formant frequency trends have been
127 reported across languages, with greater shifts for F1 than F2 or F3 (see e.g., Eklund and
128 Traunmüller, 1997; Heeren, 2015; Jovičić and Šarić, 2008).

129 Due to the described acoustic features of whispered speech, listeners have been found to be
 130 less accurate when identifying linguistic information (Konno, 2016) and emotion (Frühholz
 131 *et al.*, 2016). Whispered speech has also been found to severely degrade speaker recognition
 132 systems, which are primarily based on neutral mode speech-processing algorithms (Zhang
 133 and Hansen, 2018). However, despite lacking a fundamental frequency (F_0), whispered
 134 speech has a clearly perceivable prosodic structure. Žygis *et al.* showed that the spectral
 135 properties of consonants change during whispering to convey intonation patterns, compen-
 136 sating for the absence of a fundamental frequency. Jovičić and Šarić report longer durations
 137 for consonants. Such modifications provide evidence for cue-trading relations, where one
 138 dominant cue is substituted by the integration of multiple others, which would otherwise be
 139 less prominent when considered in isolation (Žygis *et al.*, 2017).

140 Due to its acoustical profile with a decreased amplitude envelope and reduced periodic
 141 spectro-temporal fine structure, whispered speech as compared with voiced speech should
 142 produce either similar or less interference with serial-order processing. Previous studies have
 143 shown that modulations of the spectral detail of irrelevant speech (i.e., presenting noise
 144 vocoded speech varying in the number of independently amplitude-modulated frequency
 145 bands) influences the degree of distraction, with reduced spectral fidelity (decreasing num-
 146 ber of frequency bands) attenuating disruption of serial recall (Dorsi *et al.*, 2018; Ellermeier
 147 *et al.*, 2015). Similarly, manipulations of speech prosody (e.g., emotional speech or urgent
 148 intonations) were found to increase disruption of serial recall performance (Kattner and
 149 Ellermeier, 2018; Ljungberg *et al.*, 2012), suggesting that enhanced amplitude (and fre-
 150 quency) modulations in speech intonations may increase interference with order processing

151 (note that emotional speech prosody did not affect performance on the missing-item task,
152 which does not required the retention of serial order). It has also been found that disrup-
153 tion of serial recall is determined largely by changes in vowels rather than consonants (i.e.,
154 consonant-vowel-consonant syllables are more disruptive when all components or only the
155 vowels change, compared to when a consonant changes; [Hughes et al., 2005a](#)). Hence, in par-
156 ticular due to the lower amplitude of whispered vowels ([Ito et al., 2005](#)), it could be predicted
157 that the changing-state effect on serial recall may be reduced with whispered compared to
158 voiced speech (i.e., there should be an interaction between state and phonation, see Ta-
159 ble I). More specifically, due to the lower amplitude modulations, recall accuracy should be
160 higher in the whispered changing-state condition than in the voiced changing-state condition,
161 whereas less phonation-related differences should be observed in the steady-state conditions.
162 However, there are currently no studies showing that a decrease in the depth of amplitude
163 modulations decreases distraction, and some studies found that serial recall is insensitive to
164 the overall level and intensity changes of irrelevant speech ([Ellermeier and Hellbrück, 1998](#);
165 [Tremblay and Jones, 1999](#)). In contrast, more recent findings suggest that both steady-state
166 (repeated words) and changing-state (varying words) sequences of high-intensity sound (75
167 dB(A)) are more disruptive than low-intensity sound sequences (45 dB(A)) in a serial recall
168 task ([Alikadic and Röer, 2022](#)). In line with this finding, it could be argued that due to
169 the lower amplitude modulations and overall loudness, both steady- and changing-state se-
170 quences of whispered words should be less disruptive than their voiced counterparts. In the
171 present study this was controlled partially by normalizing the amplitudes of whispered and
172 voiced speech recordings, but also by testing level and loudness as predictors of serial recall

173 accuracy. More precisely, in order to test the contribution of (psycho)acoustical properties
174 of irrelevant sound (Ellermeier and Zimmer, 2014), a regression analysis was conducted to
175 predict serial recall accuracy based on multiple signal metrics including loudness, fluctuation
176 strength, and tonality.

177 On the other hand, it could also be argued that whispered phonation increases the per-
178 ceptual demands (due to a lack of f0 cues and altered spectral fidelity) to process the speech
179 signal and to achieve speech recognition, thus reducing the resources available to process
180 the focal serial recall task. According to current speech processing models (Pichora-Fuller
181 *et al.*, 2016; Rönnberg *et al.*, 2021, 2013; Wingfield, 2016, e.g., ‘framework for understanding
182 effortful listening’ and ‘ease of language understanding’ models;) the degradation of signal
183 clarity due to the absence of f0 cues and formant alternations in whispered speech (reducing
184 speech quality and intelligibility) should impose additional cognitive processing load on pas-
185 sive listeners (e.g., speech decoding and lexical access), compared to clearly intelligible voiced
186 speech. However, this additional load is expected only in listeners who are familiar with the
187 language, because extra processing resources or ‘listening effort’ would be dedicated to irrel-
188 evant speech only when there is some degree of mismatch between the degraded (whispered)
189 speech signal and phonological representations in the listeners’ mental lexicon. Disruption
190 in the serial recall task would thus be the consequence of the enhanced listening effort re-
191 quired to process acoustically degraded, whispered speech in a comprehensible language.
192 Nevertheless, it seems more difficult to explain other effects of ‘degraded’ irrelevant speech
193 in terms of enhanced listening effort, because often degraded speech and lower speech intel-
194 ligibility results in less disruption of serial recall compared to more intelligible speech (e.g.,

195 noise-vocoded and locally time-reversed speech, [Ellermeier and Hellbrück, 1998](#); [Ellermeier](#)
196 [et al., 2015](#); [Ueda et al., 2019](#)).

197 Similar predictions could be derived from an attentional capture account though, assum-
198 ing that attention is directed to certain semantic properties or social functions associated
199 with whispered speech. As a universal, paralinguistic phenomenon found across cultures and
200 unique to human species, whispering has important social functions. These functions vary,
201 depending on whether whispering is used privately or in the public domain and influence
202 the way it is perceived ([Cirillo, 2004](#); [Cirillo and Tödt, 2005](#)). While it can be positively
203 connotated as an expression of affection in the private domain, it may elicit negative judge-
204 ments when used in the public domain. One possible explanation for this is that whispering
205 is often used to signal secrecy and confidentiality ([Laver, 1994](#)), thereby inducing mistrust
206 and social segregation and diverting the attention of non-addressees by increasing auditory
207 vigilance (compare in-group and vigilance hypothesis formulated by [Cirillo and Tödt, 2005](#)).
208 This may lead to greater attentional capture, either because whispered speech is considered
209 to be more relevant to the individual (potential self-relevance or goal-relevance of whispered
210 content), due to the greater listening effort required to process the meaning of acoustically
211 degraded (and potentially interesting) whispered background speech, or because of its dis-
212 tinctiveness in relation to the surrounding stimuli, as expressed by the salience hypothesis
213 ([Günther et al., 2017](#)).

214 Importantly, such an attentional capture mechanism should be independent of, and ad-
215 ditive to, the interference with serial-order processing produced by changing-state speech
216 (i.e., there should be no interaction between a disruptive effect of whispered speech and

217 the changing-state effect; [Hughes *et al.*, 2005b](#)). That is, whispered speech is expected to
218 cause more disruption than voiced speech both with steady- and changing-state sequences
219 of irrelevant speech (see Table I). Moreover, if the disruptive effect of whispered speech was
220 due to enhanced listening effort or the individuals' motivation to process semantic content of
221 unattended speech, then whispered speech should be more disruptive only when whispered
222 in a language that is comprehensible to the individual.

223 To the best of our knowledge, whispered speech has only been used in one previous study
224 testing the effect of background sound on the recall of short spoken lectures, but in this study
225 whispered speech was not contrasted with loud/voiced speech (and whispering was not a
226 reliable predictor of lecture recall accuracy, [Zeamer and Fox Tree, 2013](#), Exp. 3). In the
227 present series of experiments the effect of task-irrelevant whispered speech in serial recall was
228 contrasted both with normally-phonated modal speech and with a silent control condition.
229 In addition, the effect of whispering was tested both with steady-state and changing-state
230 speech. In this context, “steady state” is used as a term to describe the repetition of single
231 auditory tokens (e.g., monosyllabic words), resulting in a “steady” stream of sounds, while
232 “changing state” refers to altering auditory tokens as contained in spoken sentences. This
233 allows a test of whether whispering either (a) reduces the changing-state effect because the
234 reduced frequency and amplitude modulations interfere less with seriation, or (b) produces
235 process-independent distraction due to a diversion of attentional (e.g., triggered by social
236 functions of whispered speech) or additional cognitive demands (enhanced listening effort
237 to process whispered speech). That is, according to an interference-by-process account, the
238 changing-state effect should be reduced with whispered speech, whereas according to an

Account	Prediction	Mechanism
Interference-by-Process	<i>Phonation</i> \times <i>State Interaction</i> : Higher recall accuracy with whispered compared to voiced changing-state speech, smaller difference between whispered and voiced steady-state speech	Reduced amplitude envelope of whispered speech should provide less order information and thus cause less interference with seriation (i.e., whispered speech becomes more like a ‘steady-state’ signal)
Attentional Capture	<i>Independent main effects of Phonation and State in comprehensible language</i> : Lower recall accuracy with whispered speech compared to voiced speech and lower recall accuracy with changing-state compared to steady-state speech; <i>Main effect of State</i> , but no main effect of <i>Phonation</i> with irrelevant speech in <i>foreign language</i>	Whispered speech should cause attentional disruption due to its potential interest or enhanced listening effort, which is independent of the automatic interference produced by changing-state speech.

TABLE I. Summary of the main theoretical accounts and their predictions tested in this study and a description of the assumed mechanisms.

239 attentional capture or listening effort account there should be additive disruptive effects of
 240 changing-state and whispered speech.

241 II. EXPERIMENT 1

242 The aim of Experiment 1 was to test whether whispered speech (a) reduces the changing-
 243 state effect on serial recall as predicted by an account that assumes interference-by-process
 244 (e.g., due to reduced spectro-temporal variation / fine structure caused by the absence of fun-
 245 damental frequency cues and lower vowel amplitude) or (b) is more disruptive than voiced

246 speech as predicted by an account that predicts attentional capture by specific features
247 of whispered speech (e.g., its semantic properties). More specifically, an interference-by-
248 process account predicts an interaction between phonation and state of irrelevant speech,
249 with whispered changing-state speech being less disruptive in serial recall compared to voiced
250 changing-state speech, whereas phonation should matter less for steady-state sequences (see
251 Table I). In contrast, an attentional capture or duplex-mechanism account predicts indepen-
252 dent main effects of the phonation and state of irrelevant speech, assuming that whispered
253 speech should cause additional attentional capture and be more disruptive than voiced speech
254 regardless of whether speech is presented in steady-state or changing-state sequences. The
255 changing-state effect would thus be independent of a disruptive effect of whispered speech,
256 in particular when the language or whispered speech is comprehensible to the listener (see
257 also Table I).

258 **A. Method**

259 **1. Participants**

260 Ninety-four participants (62, female, 31 male, 1 other) were recruited at the Health and
261 Medical University campus in Potsdam, Germany. Ages ranged between 18 and 61 years
262 ($M = 24.0$, $SD = 8.6$). Participants were native speakers of German and all reported normal
263 hearing and normal or corrected-to-normal vision. The study has been conducted strictly in
264 accordance with the Ethical Principles of the Acoustical Society of America for Research.
265 All participants gave written informed consent before starting the tasks, acknowledging that

266 participation is voluntary and they were free to withdraw from the study at any time without
267 negative consequences. Participants were also informed about the scientific purpose (mainly
268 during debriefing), the potential discomfort during the task (e.g., due to cognitive demand),
269 the absence of risks to mental and physical well-being, and the confidentiality of personal
270 data. Student participants majoring in psychology ($n = 70$) were compensated with course
271 credits. Non-student participants received no compensation.

272 2. *Stimuli*

273 Two native lay speakers were recruited to record twenty unique German sentences in a
274 male and a female voice using a Behringer B1 Bundle microphone and a FMR Audio RNP
275 8380 preamplifier. The sentences were adapted from previous studies (Hughes and Marsh,
276 2020; Kattner *et al.*, 2022; Röer *et al.*, 2015) and comprised various categories such as
277 weather forecasts, traffic reports, cooking recipes, poems, operating manuals, and scientific
278 descriptions. Each sentence was spoken once with voiced phonation (normal speech) and
279 once with whispered phonation by each speaker. The speakers were instructed to adjust
280 their rate of speaking to reach about 8 s duration. The recordings were sampled at 44.1
281 kHz (16 bits). For each of the twenty (changing-state) sentences, a unique 8-s steady-
282 state sequence was created (with voiced and whispered phonation) by selecting a single
283 monosyllabic word from the sentence (e.g., ‘Hand’, 300-500 ms duration), and concatenating
284 it eight times at a rate of one word per second (we note that this creates short gaps of silence
285 between successive utterances). Thus, in total 160 sound files were created from the twenty
286 sentences (male/female speaker \times voiced/whispered phonation \times changing-/steady-state).

287 The amplitudes of all recordings were normalised in Audacity (<https://www.audacity.de/>) to
288 minimize level differences between sound conditions. Ten different sentences and steady-state
289 sequences were selected for the voiced and whispered phonation conditions, and participants
290 were presented either with the male or the female voice only. That is, forty unique speech
291 recordings were presented to each participant.

292 Exemplary FFT spectra of whispered and voiced speech are illustrated in Fig. 1. To esti-
293 mate the overall speech intensity, the A-weighted, equivalent continuous sound pressure level
294 (LA_{eq}) in dB(A) was determined for each sound file, using ArtemiS SUITE (HEAD Acoustics
295 GmbH, Herzogenrath, Germany). In addition, the psychoacoustic metrics ‘Zwicker’ loudness
296 (cubic average) as per DIN45631/A1 ([DIN Deutsches Institut für Normung e.V., 2010](#)), fluc-
297 tuation strength ([Fastl, 1982](#); [Fastl and Zwicker, 2007](#)) and sharpness as per DIN 45692 ([DIN](#)
298 [Deutsches Institut für Normung e.V., 2009](#)) were computed for each sound file. Roughness
299 and tonality were calculated according to the ECMA-418-2 (2nd) standard ([ECMA Inter-](#)
300 [national, 2022](#)). A free sound field was assumed for the calculation of all metrics. Loudness
301 reflects how loud a sound is perceived by human listeners. In contrast, decibels (dB) quan-
302 tify the physical intensity of a sound. Sharpness is another perceptual attribute, related
303 to the spectral content of sounds and in particular the high frequency components. Us-
304 ing the Relative Approach Method (RAM) ([Genuit, 1996](#)), the spectro-temporal changes
305 in the signal were quantified by extrapolating the signal history. Fluctuation strength and
306 roughness both reflect the sensation caused by variations in the amplitude and frequency
307 of sounds. While fluctuation strength mirrors the slow, rhythmical variations, roughness

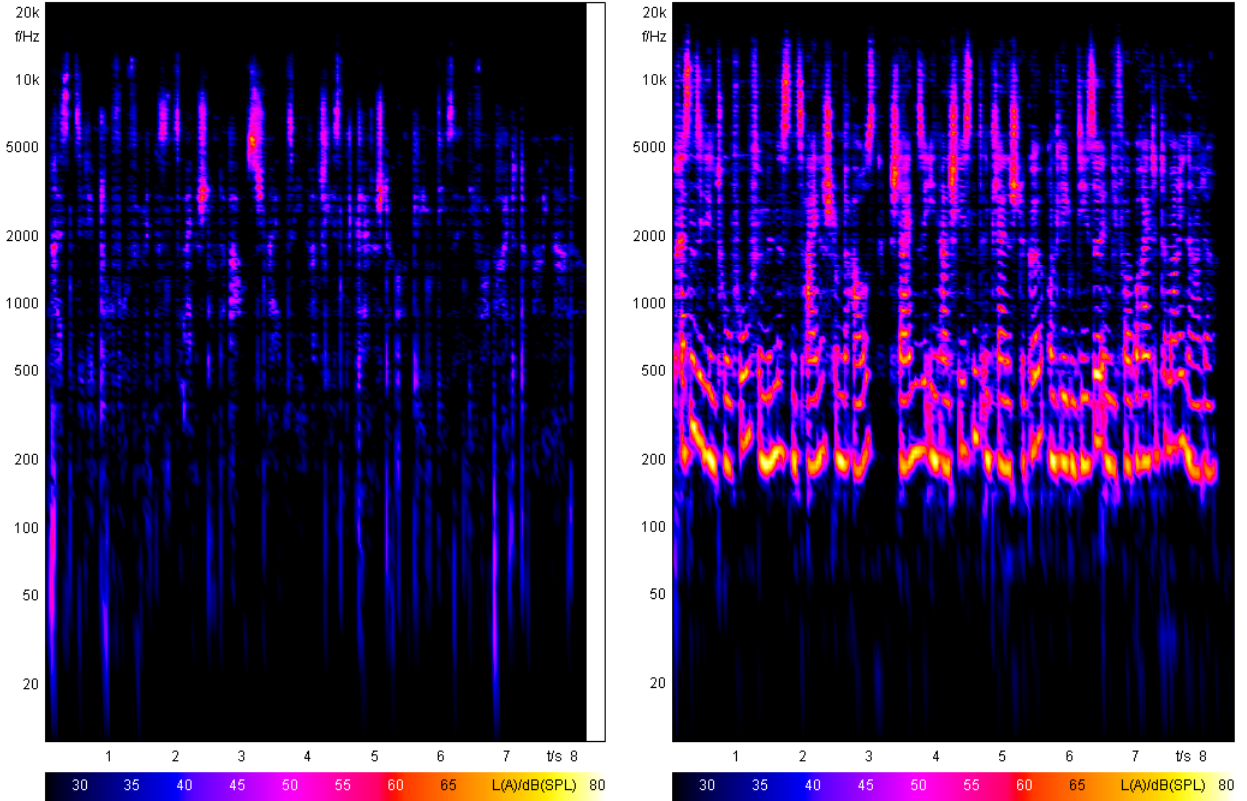


FIG. 1. FFT spectra vs. time of a changing-state sentence spoken by a female speaker with whispered (left) and voiced (right) phonation.

308 describes rapid, irregular variations. Finally, tonality indicates the relative prominence of
 309 the tonal elements within a specific noise spectrum.

310 Descriptive statistics of the psychoacoustic metrics are shown in Table II. Due to viola-
 311 tions of homogeneity of the covariance matrices ($\chi^2(198) = 768.94; p < .001$) and deviations
 312 from multivariate normality ($W = 0.87; p < .001$), non-parametric Kruskal-Wallis rank-sum
 313 tests were conducted to test for differences in psychoacoustical metrics between experimental
 314 conditions.

315 *Level* was significantly higher in voiced compared to whispered speech, $W(1) = 1716$,
 316 $p < .001$. Moreover, there was a significant level difference between steady- and changing-

317 state speech, $W(1) = 53.18$, $p < .001$ (higher levels in steady-state), and also between male
 318 and female voices, $W(1) = 4.43$, $p = .035$.

319 *'Zwicker' loudness* was significantly higher for voiced than for whispered speech, $W(1) =$
 320 49.57 , $p < .001$, and steady-state sequences are louder than changing-state sentences,
 321 $W(1) = 33.05$, $p < .001$. The speaker difference in loudness was not significant, $W(1) = 1.30$,
 322 $p = .254$. We note that the loudness differences between state and voice conditions are
 323 rather small in magnitude and any detrimental effect of loudness on serial recall would work
 324 against the main hypotheses that the softer whispered speech and changing-state speech will
 325 be more disruptive than voiced speech and steady-state speech. Moreover, reducing loudness
 326 differences through normalization could have removed the characteristic attention-capturing
 327 properties of whispered speech.

328 *Sharpness* was significantly higher for whispered than for voiced speech, $W(1) = 66.72$,
 329 $p < .001$, reflecting the larger amount of high-frequency energy in whispered speech. There
 330 was also a speaker difference, $W(1) = 21.35$, $p < .001$, but no difference between steady-
 331 and changing-state sequences in sharpness, $W(1) = 0.60$, $p = .441$.

332 The difference in *roughness* between whispered and voiced speech was also significant,
 333 , $W(1) = 14.17$, $p < .001$, likely attributed to the absence of an F_0 in whispered speech.
 334 Roughness was also higher in male than in female speech, , $W(1) = 56.50$, $p < .001$, but there
 335 was no roughness difference between steady- and changing-state sequences, $W(1) = 0.20$,
 336 $p = .656$.

337 *Fluctuation strength* in turn was significantly higher in voiced speech compared to whis-
 338 pered speech, $W(1) = 47.22$, $p < .001$, indicating a reduced amplitude envelope with whis-

339 pered phonation, and steady-state speech was significantly more fluctuating than changing-
340 state speech, $W(1) = 89.55$, $p < .001$ (as to be expected due to the silent gaps between
341 successive words in steady-state streams). There was no significant speaker difference in
342 fluctuation strength though, $W(1) = 1.02$, $p = .31$.

343 There was also a significant difference in the spectro-temporal variation quantified via
344 the *Relative Approach* metric (an extrapolation method) between steady- and changing-
345 state sequences, $W(1) = 20.91$, $p < .001$, as well as between whispered and voiced speech,
346 $W(1) = 30.68$, $p < .001$, reflecting more spectral and temporal variation in voiced speech.
347 There was no speaker difference in spectro-temporal variation, $W(1) = 1.98$, $p = .159$.

348 Finally, *tonality* was higher in voiced than in whispered speech, $W(1) = 79.88$, $p < .001$,
349 as well as in steady-state speech compared to changing-state sentences, $W(1) = 4.41$, $p =$
350 $.036$. These distinct differences in vocal quality are also visible in Figure 1. Tonality was
351 also higher in female speech than in male speech, $W(1) = 9.75$, $p = .002$.

352 **3. Apparatus**

353 The study was conducted in a single-walled sound-attenuated listening booth (Studiobox
354 GmbH, Munich, Germany). The experiment ran on a Lenovo Thinkstation P350 desktop
355 computer and the experimental routines were programmed in Python utilizing the PsychoPy
356 package (Peirce *et al.*, 2019). Visual stimuli were presented on a BenQ GW2780 IPS screen
357 (27 in).

358 Sounds were D/A converted by an ESI MAYA44 eX PCIe sound card (ESI Audiotechnik,
359 Leonberg, Germany) passed through a Behringer Powerplay HA8000 amplifier (Behringer,

TABLE II. Mean psychoacoustic metrics of the 20 changing-state sentences and 20 monosyllabic steady-state word sequences, each spoken aloud (voiced) and whispered by a male and female speaker (standard deviations in parentheses).

Parameter	Speaker	Changing-State		Steady-State	
		whispered	voiced	whispered	voiced
LA _{eq} (dB(A))	male	71.62 (2.47)	74.79 (1.63)	78.16 (2.89)	79.60 (2.89)
	female	71.28 (2.51)	75.52 (1.22)	74.36 (3.32)	77.18 (2.30)
Loudness (sone)	male	17.63 (2.31)	24.75 (2.53)	28.07 (5.08)	30.52 (4.99)
	female	16.64 (2.50)	27.98 (2.65)	22.07 (4.83)	28.77 (4.50)
Sharpness (acum)	male	1.66 (0.11)	1.34 (0.09)	1.68 (0.16)	1.47 (0.22)
	female	1.89 (0.10)	1.54 (0.09)	1.78 (0.14)	1.59 (0.16)
Roughness (asper)	male	0.37 (0.08)	0.88 (0.16)	0.33 (0.08)	1.05 (0.39)
	female	0.31 (0.04)	0.26 (0.03)	0.28 (0.10)	0.35 (0.09)
Fluctuation Strength (vacil)	male	0.20 (0.04)	0.55 (0.15)	0.56 (0.16)	1.34 (0.24)
	female	0.18 (0.03)	0.36 (0.07)	0.62 (0.19)	1.32 (0.21)
Rel. Approach (cPa)	male	38.72 (5.24)	54.58 (6.31)	39.24 (7.33)	43.05 (8.31)
	female	40.77 (4.40)	48.69 (5.34)	37.02 (4.56)	38.94 (6.34)
Tonality (<i>t.u.</i> _{HMS})	male	0.21 (0.05)	0.39 (0.12)	0.34 (0.09)	0.41 (0.15)
	female	0.25 (0.04)	0.91 (0.15)	0.27 (0.11)	1.12 (0.41)

360 Penang, Malaysia) and played diotically via Beyerdynamic DT 990 PRO headphones (Beyer-
361 dynamic, Heilbronn, Germany) at an overall average playback level of 70 dB(A), (with inter-
362 sentence variability accounting for the differences in state and phonation, see Table II). This
363 playback level deviated slightly from the original sound pressure level of the recordings, but
364 was deemed comfortable for participants.

4. *Experimental Design and Procedure*

A 2 (State: steady, changing) \times 2 (Phonation: voiced, whispered) experimental design was implemented. Silence was presented as a control condition to assess possible disruptive effects of steady-state sound (Bell *et al.*, 2019). There were ten repetitions of each auditory condition, resulting in a total of 50 trials that were presented in fully randomized order. Half of the participants ($n = 47$) were presented with irrelevant speech in the male voice, and the other half was presented with the female voice only.

Participants were instructed to memorize the order of eight consonants presented on the screen while ignoring the sound that was played via headphones. Participants started each trial at their own pace by pressing the space bar. Then an empty white square was presented in the center of the black screen for 1 s before the eight to-be-remembered consonants were presented successively within the square. The consonants were drawn randomly without replacement from ‘F’, ‘G’, ‘K’, ‘L’, ‘M’, ‘P’, ‘Q’, ‘S’, and ‘T’. Each consonant was presented for 800 ms and followed by a 200-ms inter-stimulus interval showing the empty square. Irrelevant sound was presented during the visual presentation of consonants (8 s). After a silent retention interval of 6 s (showing a blank screen), a 3 \times 3 response matrix was presented on the screen showing all nine consonants arranged alphabetically. Participants were prompted to click the consonants in the memorized order. The sequence of clicked consonants was presented on the screen (above the matrix). Participants were able to click consonants multiple times, but they could not correct their previous responses. After the last click response, the number of consonants that were recalled in the correct serial position

386 was presented as visual feedback for 1.5 s (e.g., ‘Trial 3: 6 correct’). The next trial started
 387 immediately after the feedback. After the 10th, 20th, 30th, and 40th trial, an additional text
 388 prompt was presented on the screen, indicating that participants could now take a short
 389 break before proceeding with the next trial.

390 B. Results

391 The serial recall accuracy in the five auditory conditions is illustrated in Fig. 2, both
 392 averaged and across serial positions. A 2 (state: steady, changing) \times 2 (phonation: voiced,
 393 whispered) \times 8 (serial position: 1-8) repeated-measures ANOVA on recall accuracy revealed
 394 a significant main effect of state, $F(1, 93) = 50.13$, $MSE = 0.05$, $p < .001$, $\hat{\eta}_p^2 = .350$
 395 (i.e., a changing-state effect: lower recall accuracy with changing-state speech compared to
 396 steady-state words), and a significant main effect of phonation, $F(1, 93) = 6.27$, $MSE = 0.05$,
 397 $p = .014$, $\hat{\eta}_p^2 = .063$, with lower recall accuracy during whispered speech ($M = .53$, $SD = .12$)
 398 than during voiced speech ($M = .55$, $SD = .12$). The relative decrement in performance com-
 399 pared to silence ($(accuracy_{silence} - accuracy_{speech}) / accuracy_{silence}$, cf. [Ellermeier and Zimmer,](#)
 400 [2014](#)) was 12.9% for whispered speech, compared to 9.6% for voiced speech. There was no
 401 interaction between state and phonation, $F(1, 93) = 0.17$, $MSE = 0.04$, $p = .680$, $\hat{\eta}_p^2 = .002$
 402 As can be seen in Fig. 2A, whispered speech produced similar disruption of serial recall
 403 compared to voiced speech, regardless of whether speech consisted of steady-state words
 404 and changing-state sentences.

405 As to be expected, the ANOVA also revealed a significant serial position effect, $F(3.42, 317.70) =$
 406 317.75 , $MSE = 0.10$, $p < .001$, $\hat{\eta}_p^2 = .774$, with higher accuracy for items from the beginning

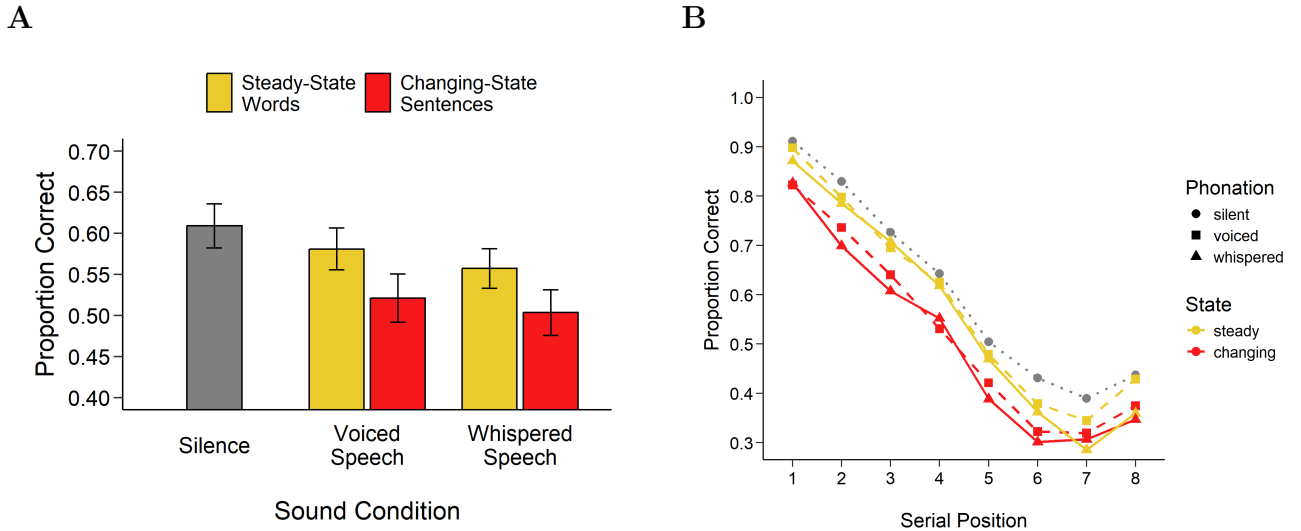


FIG. 2. (A) Mean serial recall accuracy in silence and when either steady-state or changing-state speech was presented with voiced or whispered phonation during item encoding in Experiment 1. Error bars represent 95% confidence intervals. (B) Serial recall accuracy in the five irrelevant sound conditions as a function of serial position.

407 of the list (position 1: $M = 0.85$ [0.83, 0.88], position 2: $M = 0.75$ [0.73, 0.78], position 3:
 408 $M = 0.66$ [0.63, 0.69], position 4: $M = 0.58$ [0.55, 0.62], position 5: $M = 0.44$ [0.40, 0.48],
 409 position 6: $M = 0.34$ [0.31, 0.38], position 7: $M = 0.31$ [0.28, 0.35]; 95% CIs in brackets)
 410 as well as a small recency effect (position 8: $M = 0.38$ [0.34, 0.41]). It was further tested
 411 whether the effect of whispering and changing-state sound differs for items in different serial
 412 positions (compare Fig. 2B). To that effect, the ANOVA revealed a significant interaction
 413 between state and serial position, $F(5.31, 493.77) = 4.14$, $MSE = 0.02$, $p < .001$, $\hat{\eta}_p^2 = .043$,
 414 but not between phonation and serial position, $F(5.12, 476.25) = 1.70$, $MSE = 0.02$,
 415 $p = .130$, $\hat{\eta}_p^2 = .018$. A planned contrasts analysis corrected for multiple comparisons
 416 (Benjamini and Hochberg, 1995) revealed that the changing-state effect was significant at
 417 serial positions 1 to 6 ($p_{\text{BH}(8)} < .001$) and 8 ($p_{\text{BH}(8)} = .017$), but not at serial position 7

418 ($p_{\text{BH}(8)} = .886$). There was no three-way interaction between state, phonation, and serial
 419 position, $F(6.12, 569.19) = 1.77$, $MSE = 0.02$, $p = .101$, $\hat{\eta}_p^2 = .019$.

420 C. Discussion

421 Experiment 1 demonstrated that whispered speech produced more disruption in a serial
 422 recall task compared to speech presented with (louder) voiced phonation. As expected,
 423 changing-state speech (full German sentences) was also more disruptive than steady-state
 424 speech consisting of repetitions of a single monosyllabic German word. Interestingly, these
 425 two effects seem to be independent, as indicated by the absence of an interaction. While the
 426 changing-state effect is most likely due to interference between the order information in the
 427 auditory stream and deliberate serial-order processing, the “whispering effect” may be due
 428 to attentional capture elicited either by the potential meaning of whispered information or
 429 the enhanced listening effort required to process the semantic content of whispered speech
 430 in a comprehensible language (as predicted by speech processing accounts, [Pichora-Fuller](#)
 431 [et al., 2016](#); [Rönnberg et al., 2021, 2013](#); [Wingfield, 2016](#)). To test this last assumption,
 432 a second experiment was conducted in which we tried to replicate the disruptive effect of
 433 whispered speech that is presented in a language that is foreign to the listener, making it
 434 incomprehensible. If the attentional disruption was due to enhanced listening effort in case of
 435 acoustically degraded but comprehensible whispered speech, then it should disappear when
 436 participants perceive the language as an incomprehensible, foreign language, because in this
 437 case there would be no mismatch between the task-irrelevant speech signal and phonological
 438 representations stored in the mental lexicon.

439 **III. EXPERIMENT 2**

440 Experiment 2 was a close replication of Experiment 1, but with a sample of participants,
441 who did not understand the irrelevant speech language (German).

442 **1. Participants**

443 A power analysis based on the effect size for the whispering effect observed in Experiment
444 1 ($\hat{\eta}_p^2 = .063$) revealed that a sample size of $N = 51$ is required to reach a statistical
445 power of $1 - \beta = .95$ ($\alpha = .05$) for the detection of a two-level main effect in a repeated-
446 measures ANOVA. Fifty-one participants (42 women) who did not speak or understand
447 German were recruited either at the University of Lincoln, UK ($n = 44$), or at Ludwig
448 Maximilian Universität München, Germany ($n = 7$). Ages ranged between 18 and 53 years
449 ($M = 29.3$; $SD = 11.8$). All participants reported normal hearing and normal or corrected-
450 to-normal vision. Most participants of Experiment 2 were native speakers of English, but
451 there were also a few native speakers of other languages (e.g., Chinese and Spanish). We
452 also note that an additional data analysis including only the subsample of native speakers of
453 English – not including speakers of a logographic language such as Chinese – produced the
454 same overall pattern of results. All participants confirmed not speaking or understanding
455 the German language. The study has received ethics approval by the ethics committee of
456 the University of Lincoln (ref: 33415). All participants gave written informed consent before
457 starting the task. Participants of Experiment 2 were compensated with course credit.

2. *Stimuli and Apparatus*

The set of German speech recordings from Experiment 1 was used also for Experiment 2, but 16 unique changing-state and steady-state recordings were selected. Half of the speech samples were presented with voiced phonation and half were presented with whispered phonation. Each sentence or word sequence was selected once in the male and once in the female voice, thus generating 16 unique speech recordings for each auditory condition (state \times phonation).

The experiment was conducted on an HP EliteDesk computer in a testing cubical at the University of Lincoln. Visual stimuli were presented on an HP EliteDisplay E240 screen (24 in). An Intel Realtek audio controller was used and sounds were played dichotically via Sony MDR-ZX110 headphones at a level similar to Experiment 1 (approximately 70 dB(A) on average). The experiment was programmed in Python using PsychoPy ([Peirce et al., 2019](#)).

3. *Design and Procedure*

The experimental design was the same as in Experiment 1, using five different auditory conditions (silence, voiced/whispered steady-state words, voiced/whispered changing-state speech). The procedure was identical to Experiment 1, except that the number of repetitions per experimental condition was increased to 16, resulting in a total of 80 trials. As in Experiment 1, unique sentences or unique steady-state words were presented on each trial. Moreover, half of the speech trials were presented by the male and female voice, respectively. The trial structure was also identical to Experiment 1, except that after each trial partici-

478 pants were asked to give a confidence judgment by clicking on a scale from 0 to 8, indicating
 479 “how many letters they thought to have recalled in the correct position”. Feedback on the
 480 actual number of correct letters was presented after the confidence judgment.

481 A. Results

482 1. *Whispering and changing-state effects with foreign language*

483 In contrast to Experiment 1, an equivalent 2 (state: steady, changing) \times 2 (phonation:
 484 voiced, whispered) \times 8 (serial position) repeated-measures ANOVA revealed a significant
 485 main effect of state, $F(1, 51) = 7.37$, $MSE = 0.03$, $p = .009$, $\hat{\eta}_p^2 = .126$, but no significant
 486 main effect of phonation, $F(1, 51) = 2.85$, $MSE = 0.03$, $p = .097$, $\hat{\eta}_p^2 = .053$. As can be seen
 487 in Fig. 3A, in participants to whom the language is incomprehensible, whispered German
 488 speech tended to be less disruptive ($M = .53$; $SD = .12$) compared to voiced German speech
 489 ($M = .51$, $SD = .10$). The relative decrement in performance (compared to silence) was
 490 12.8% for voiced speech and 10.4% for whispered speech. There was also no interaction
 491 between phonation and state in Experiment 2, $F(1, 51) = 0.14$, $MSE = 0.03$, $p = .713$,
 492 $\hat{\eta}_p^2 = .003$.

493 The ANOVA also revealed a significant serial position effect, $F(2.54, 129.62) = 204.73$,
 494 $MSE = 0.13$, $p < .001$, $\hat{\eta}_p^2 = .801$, as well as an interaction between state and serial
 495 position, $F(5.21, 265.50) = 2.40$, $MSE = 0.01$, $p = .035$, $\hat{\eta}_p^2 = .045$. According to a planned
 496 contrasts analysis, the changing-state effect was significant only at the early serial positions
 497 1 ($p_{\text{BH}(8)} = .046$), 2 ($p_{\text{BH}(8)} = .008$) and barely at position 3 ($p_{\text{BH}(8)} = .063$), but not at the

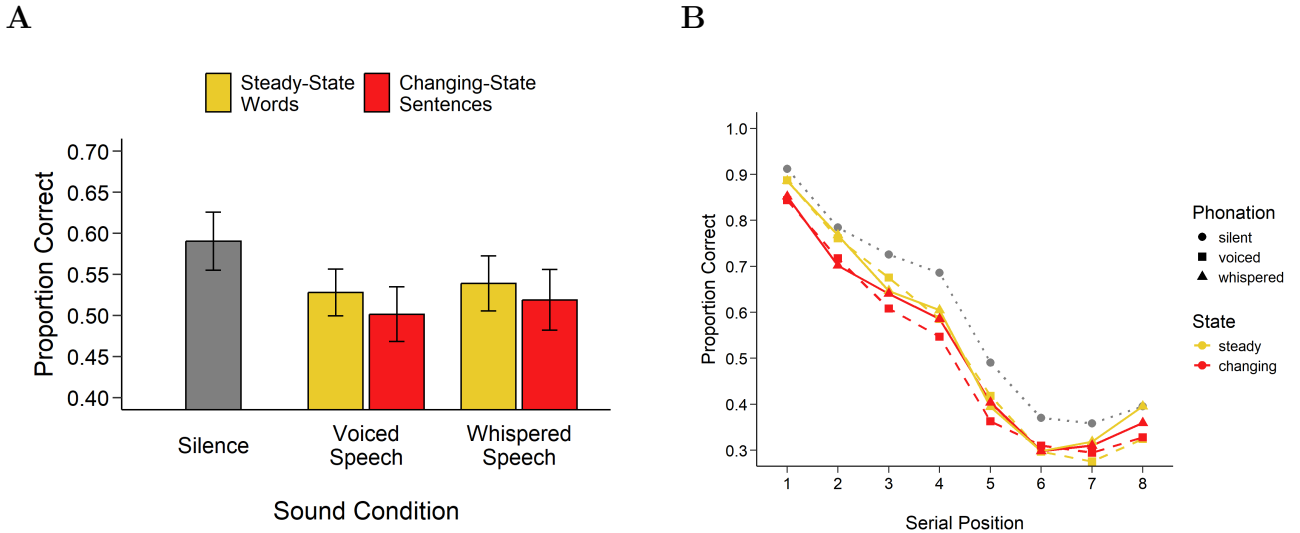


FIG. 3. **(A)** Mean serial recall accuracy in silence and when either steady-state or changing-state speech was presented with voiced or whispered phonation in a foreign language (German) during encoding in Experiment 2. Error bars represent 95% confidence intervals. **(B)** Serial recall accuracy in the five irrelevant sound conditions as a function of serial position.

498 later serial positions ($p_{\text{BH}(8)} \geq .179$). Consistent with Experiment 1, the interaction between
 499 phonation and serial position was not significant, $F(5.46, 278.29) = 2.04$, $MSE = 0.01$, $p =$
 500 $.067$, $\hat{\eta}_p^2 = .039$. However, we note that there was a non-significant trend towards whispered
 501 speech being less disruptive than voiced speech at the last serial position, $t(51) = -2.72$,
 502 $p_{\text{BH}(8)} = .071$, whereas all other contrasts were clearly non-significant ($p_{\text{BH}(8)} \geq .162$). This
 503 indicates that task-irrelevant whispered speech in a non-comprehended language is equally
 504 disruptive as voiced speech to the memorization of items from the beginning and the middle
 505 of the list, but it may restore the recency effect (compare Fig. 3B). There was also no
 506 significant three-way interaction, $F(4.84, 246.66) = 2.08$, $MSE = 0.01$, $p = .070$, $\hat{\eta}_p^2 = .039$.

507 2. *Metacognitive confidence*

508 To assess metacognitive awareness of the disruptive effects of changing-state and whis-
509 pered speech, participants of Experiment 2 were also asked to indicate their confidence after
510 each trial. A one-way repeated-measures ANOVA revealed a significant difference in confi-
511 dence judgments between the five auditory conditions, $F(3.23, 164.95) = 19.69$, $MSE = 0.16$,
512 $p < .001$, $\hat{\eta}_p^2 = .279$ (for descriptive statistics, see Table III). Planned contrasts revealed that
513 confidence was higher in silence compared to both steady-state ($p < .001$) and changing-state
514 speech ($p < .001$), but there was no significant difference in confidence between steady-state
515 and changing-state conditions ($p = .082$; note however that it would be premature to con-
516 clude that participants did not notice the difference between the two conditions, see Bell
517 *et al.*, 2022; Kattner and Bryce, 2022; Röer *et al.*, 2017b), nor between voiced and whispered
518 speech conditions ($p = .217$). This indicates that participants were aware of the general dis-
519 ruption by task-irrelevant speech, but they did not notice the stronger impairment by specific
520 types of speech (e.g., changing-state speech).

TABLE III. Means and standard deviations of confidence judgments of serial recall in silence and during the presentation of steady-state words or changing-state sentences with voiced or whispered phonation in Experiment 2.

	silence	steady-state		changing-state	
		voiced	whispered	voiced	whispered
<i>M</i>	4.07	3.61	3.66	3.52	3.59
<i>SD</i>	1.23	1.02	1.06	1.07	1.13

521 **3. Cross-experiment analysis**

522 To directly compare the effects of whispered speech in a comprehensible and incompre-
523 hensible language (i.e., Experiment 1 vs. 2), an additional 2 (experiment) \times 2 (state) \times 2
524 (phonation) mixed-factors ANOVA was conducted with experiment as a between-subjects
525 factor and state and phonation as within-subject factors. The analysis revealed that there
526 was no main effect of phonation, $F(1, 144) = 0.25$, $MSE = 0.01$, $p = .618$, $\hat{\eta}_p^2 = .002$, but a
527 significant interaction between phonation and experiment, $F(1, 144) = 7.50$, $MSE = 0.01$,
528 $p = .007$, $\hat{\eta}_p^2 = .049$, suggesting that whispered speech was more disruptive than voiced
529 speech when the language is intelligible (Experiment 1), but not when it is incomprehen-
530 sible to participants (Experiment 2). Planned contrasts (corrected according to [Benjamini](#)
531 [and Hochberg, 1995](#)) revealed that there was a significant difference between whispered and
532 voiced phonation in Experiment 1, $t(144) = 2.71$, $p_{\text{BH}(2)} = .015$, but not in Experiment
533 2, $t(144) = -1.39$, $p_{\text{BH}(2)} = .165$. Interestingly, in addition to the main effect of state,
534 $F(1, 144) = 40.79$, $MSE = 0.01$, $p < .001$, $\hat{\eta}_p^2 = .221$, the ANOVA also revealed a signifi-
535 cant interaction between state and experiment, $F(1, 144) = 7.50$, $MSE = 0.01$, $p = .007$,
536 $\hat{\eta}_p^2 = .049$, indicating that the magnitude of the changing-state effect differed also between
537 experiments. Planned contrasts revealed that the changing-state effect was significant in
538 both experiments, but larger in Experiment 1, $t(144) = 7.60$, $p_{\text{BH}(2)} < .001$, than in Experi-
539 ment 2, $t(144) = 2.31$, $p_{\text{BH}(2)} = .022$. There were no other significant effects.

540 **4. *Psychoacoustical predictors***

541 To test whether the disruption of serial recall can be predicted by the psychoacoustic
 542 properties of irrelevant speech, a backward stepwise multiple linear regression analysis was
 543 conducted to predict the average serial recall accuracy associated with each sound file that
 544 was presented in the two experiments. In addition to the three dummy-coded categor-
 545 ical predictors phonation (0 = voiced, 1 = whispered), speaker gender (0 = male, 1 =
 546 female) and experiment (0 = Exp. 1, 1 = Exp. 2), the z -transformed psychoacoustic met-
 547 rics ‘Zwicker’ loudness, sharpness, roughness, fluctuation strength, relative approach (i.e.,
 548 spectro-temporal variation determined with the ‘Relative Approach’ method) and tonality
 549 were entered as continuous predictor variables. The starting model also contained inter-
 550 action terms for each psychoacoustic metric with experiment (except loudness due to a
 551 high variable inflation factor), but other interaction terms were not included due to multi-
 552 collinearity (as indicated by variable inflation factors). The regression analysis revealed that
 553 the best-fitting model includes only the intercept ($b = 0.55$, 95% CI [0.50, 0.60]) and exper-
 554 iment ($b = 0.01$, 95% CI [-0.02, 0.04]), fluctuation strength ($\beta = 0.05$, 95% CI [0.03, 0.08])
 555 and relative approach ($\beta = -.001$, $t(139) = -1.97$, $p = .050$) as well as the interaction
 556 term between fluctuation strength and experiment ($\beta = -0.05$, 95% CI [-0.08, -0.01]) as
 557 predictors of serial recall accuracy, $R^2 = .18$, $F(4, 139) = 7.43$, $p < .001$.

558 To further investigate the predictive power of individual psychoacoustic metrics while
 559 avoiding multicollinearity, additional backward step-wise regression analyses were conducted

560 for each psychoacoustic predictor variable including the respective two-way interaction terms
 561 with experiment, phonation and speaker gender.

562 ‘*Zwicker*’ loudness was found to be a small but significant predictor, $\beta = 0.00$, 95% CI
 563 [0.00, 0.00], $t(141) = 2.76$, $p = .007$, and together with an interaction term with experiment,
 564 $\beta = 0.00$, 95% CI [0.00, 0.00], $t(141) = -2.86$, $p = .005$, it accounted for about 8% of the
 565 variance in serial recall accuracy, $R^2 = .08$, $F(2, 141) = 6.36$, $p = .002$. As can be seen
 566 in Fig. 4, increasing loudness was associated with better performance in the serial recall
 567 task, and this relationship was stronger with comprehensible speech in Experiment 1 than
 568 in Experiment 2. The same predictive relationships were found also for a model including
 569 A-weighted *sound pressure level* and its interaction with experiment, which accounted for
 570 even 10% of the variance, $R^2 = .10$, $F(2, 141) = 7.81$, $p < .001$ (see Fig. 4).

571 *Fluctuation strength* was also found to be an important predictor, with increasing fluc-
 572 tuation strength predicting higher recall accuracy (see Fig. 5; in contrast to the assumption
 573 of stronger disruption of serial recall by sounds of higher fluctuation strength, Schlittmeier
 574 *et al.*, 2012), $\beta = 0.05$, 95% CI [0.03, 0.07], $t(140) = 5.08$, $p < .001$. In addition, there was a
 575 significant interaction term between fluctuation strength and experiment, $\beta = -0.04$, 95%
 576 CI [-0.06, -0.02], $t(140) = -3.42$, $p < .001$, indicating that the predictive power of fluctu-
 577 ation strength was stronger in Experiment 1 (see Fig. 5). This suggests that the positive
 578 relationship between fluctuation strength and recall accuracy is stronger in participants who
 579 are able to understand the distractor language, and it also reflects the fact that there was
 580 only a changing-state effect but no whispering effect in the absence of speech comprehension
 581 (Experiment 2). The model also contained a small, but non-significant interaction term

582 between fluctuation strength and phonation, $\beta = 0.03$, 95% CI [0.00, 0.06], $t(140) = 1.84$,
 583 $p = .068$, but this may be biased by the lower and smaller range of fluctuation strength in
 584 whispered speech sounds (see Fig. 5). Together, the fluctuation strength model accounted
 585 for 17% of the variance in serial recall, $R^2 = .17$, $F(3, 140) = 9.55$, $p < .001$.

586 The regression model with *relative approach* as an indicator of spectro-temporal variation
 587 in the signal accounted for 8% of the variance in serial recall, $R^2 = .08$, $F(3, 140) = 3.85$,
 588 $p = .011$, with higher relative approach predicting lower serial recall accuracy (see Fig. 5).
 589 However, both the main effect of relative approach, $\beta = 0.00$, 95% CI [0.00, 0.00], $t(140) =$
 590 -2.03 , $p = .045$, and its interaction with experiment, $\beta = 0.00$, 95% CI [0.00, 0.00], $t(140) =$
 591 -2.33 , $p = .021$, were both rather small though significant predictors.

592 Psychoacoustic *roughness* was only a small and non-significant predictor of serial recall
 593 accuracy, $\beta = 0.02$, 95% CI [-0.01, 0.05], $t(141) = 1.45$, $p = .148$, but together with the
 594 interaction with experiment, $\beta = -0.04$, 95% CI [-0.07, -0.01], $t(141) = -2.94$, $p = .004$,
 595 it also accounted for some variance in serial recall accuracy, $R^2 = .06$, $F(2, 141) = 4.32$,
 596 $p = .015$. This may indicate that the rougher voiced speech sounds were less disruptive
 597 with comprehensible speech in Experiment 1, whereas they tend to be more disruptive with
 598 incomprehensible speech in Experiment 2 (see also Fig. 5).

599 The model with *sharpness*, $\beta = 0.00$, 95% CI [-0.04, 0.05], $t(141) = 0.17$, $p = .867$, and
 600 its interaction term with experiment, $\beta = -0.01$, 95% CI [-0.02, 0.00], $t(141) = -2.00$,
 601 $p = .047$, did not achieve a significant fit and accounted only for a small portion of the
 602 variance in serial recall, $R^2 = .03$, $F(2, 141) = 2.01$, $p = .138$.

603 The same is true for *tonality*, $R^2 = .03$, $F(2, 141) = 2.53$, $p = .084$, which tends to be
604 positively related with serial recall accuracy, $\beta = 0.03$, 95% CI [0.00, 0.06], $t(141) = 1.90$,
605 $p = .060$. Moreover, the relationship between tonality and recall also differed between
606 experiments, $\beta = -0.03$, 95% CI [-0.06, 0.00], $t(141) = -1.98$, $p = .049$.

607 B. Discussion

608 Experiment 2 revealed that whispered German speech was equally disruptive as voiced
609 speech to listeners who did not understand the language. This suggests that listening to
610 whispered phonation may not demand additional listening effort compared to voiced phona-
611 tion if participants cannot understand the language of the whispered speech. If anything,
612 it tends to be even less disruptive, presumably due to the reduced spectro-temporal varia-
613 tion and lower amplitude modulations, indicating that with an incomprehensible language,
614 disruption may be driven primarily by psychoacoustic properties of irrelevant sound.

615 While whispered speech was equally disruptive as voiced speech, indicating no (addi-
616 tional) attentional capture with incomprehensible speech, Experiment 2 still revealed a clear
617 changing-state effect, indicating interference-by-process (Jones and Tremblay, 2000). How-
618 ever, although almost the same speech recordings were used, the size of the changing-state
619 effect was larger in Experiment 1 than in Experiment 2, indicating that the interference
620 with order processing may be more pronounced in a language that is comprehensible to
621 participants. This could be explained with more efficient auditory grouping of speech to-
622 kens in a familiar language, thus forming a more stable auditory stream and in case of a
623 changing-state stream more disruption of serial-order processing.

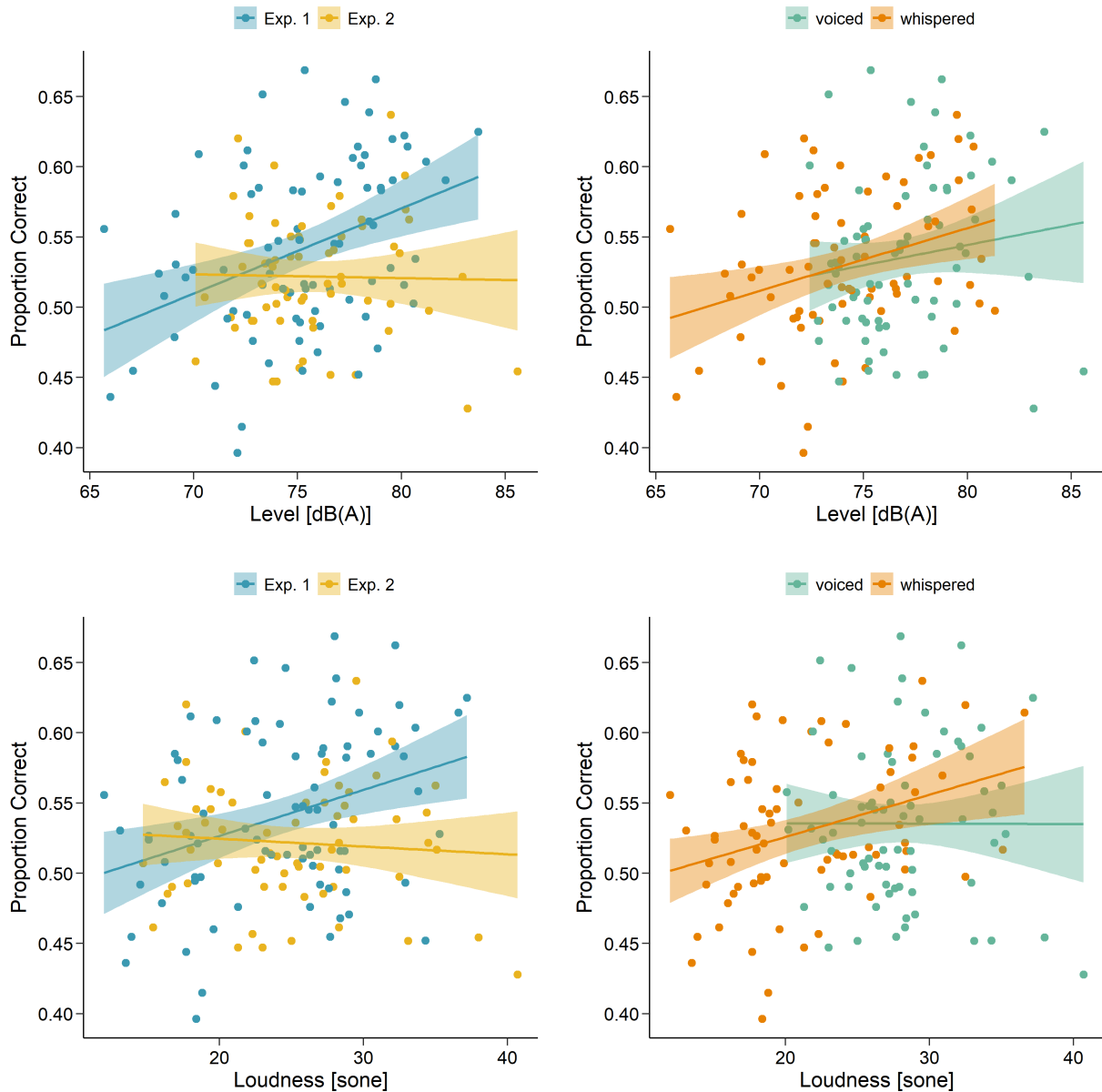


FIG. 4. Serial recall accuracy predicted by the sound pressure level [dB(A)] and ‘Zwicker’ loudness [sone] of whispered and voiced task-irrelevant German speech presented in Experiment 1 (German listeners) and 2 (foreign listeners).

624 IV. GENERAL DISCUSSION

625 The present study investigated whether task-irrelevant whispered speech produces more
 626 or less disruption of serial recall from visual-verbal short-term memory compared to voiced

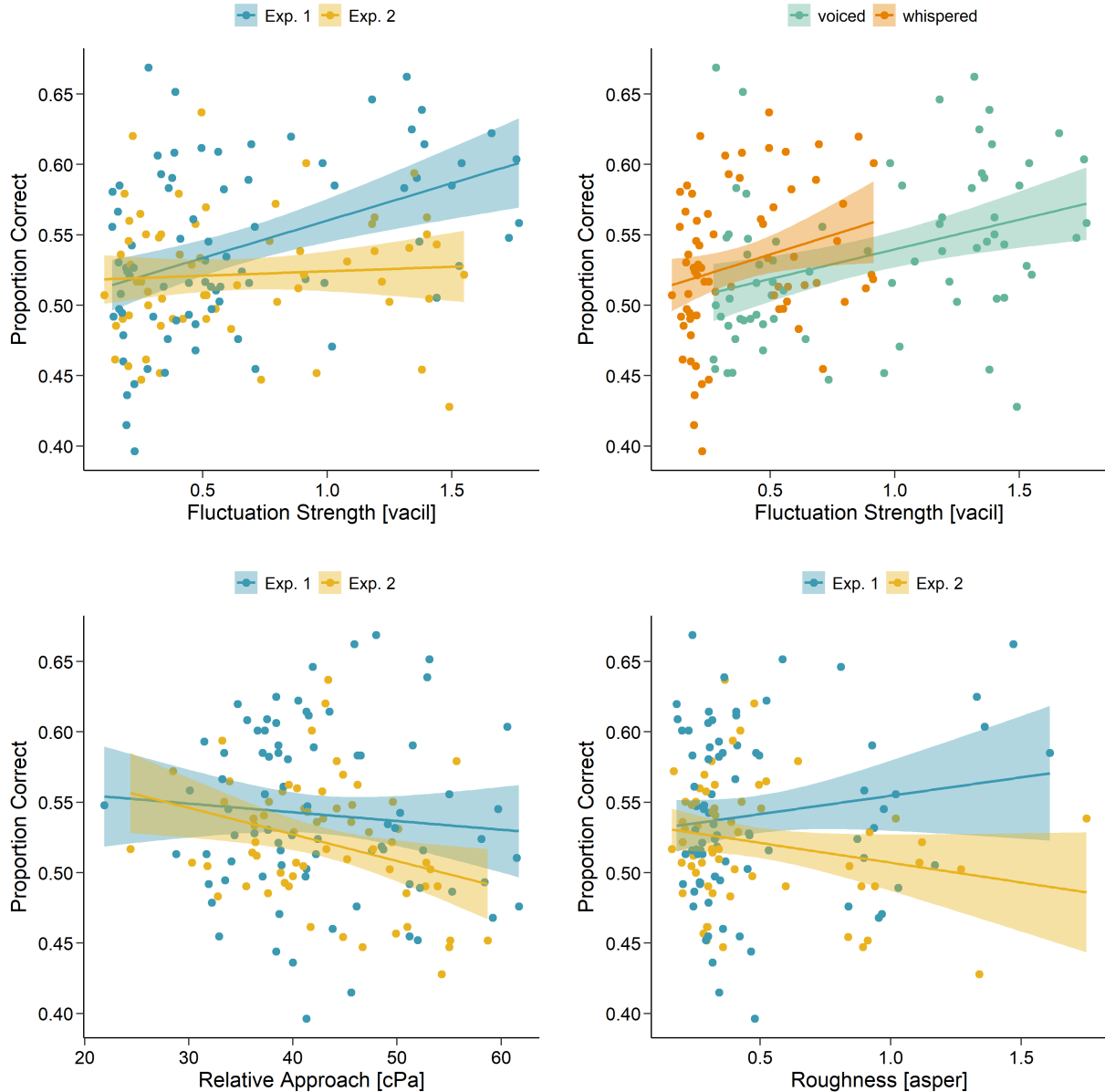


FIG. 5. Serial recall accuracy predicted by fluctuation strength [vacil], relative approach [cPa], and roughness [asper] of whispered and voiced German speech presented in Experiment 1 (German listeners) and 2 (foreign listeners).

627 speech. Specifically, according to an attentional account of auditory distraction, whispered
 628 speech could be expected to capture more attentional-cognitive resources than voiced speech,
 629 either due to its potential self-relevance (e.g., Röer *et al.*, 2013, 2017a) or because it re-

630 quires additional listening effort to process ‘degraded’ whispered speech in a comprehensi-
 631 ble language (e.g., due to the missing harmonic structure). In contrast, according to an
 632 interference-by-process account, less disruption by whispered speech would be expected,
 633 given that whispered speech sounds with a reduced level, amplitude envelope and fluctu-
 634 ation strength may cause less interference with serial-order processing (compare [Alikadic
 635 and Röer, 2022; Jones *et al.*, 2000](#)). In line with the attentional account, it was found in
 636 Experiment 1 that whispered speech in a comprehensible language causes about 35% more
 637 disruption of serial recall than voiced speech (i.e., the relative accuracy decrements were
 638 9.6% with voiced speech and 12.9% with whispered speech). Interestingly, the disruptive
 639 effect of whispered speech was found to be independent of the changing-state effect. That
 640 is, whispered speech was more disruptive both in steady-state sequences (repetitions of a
 641 single monosyllabic word) and changing-state spoken sentences, and changing-state speech
 642 was more disruptive than steady-state speech regardless of whether the phonation type was
 643 whispered or voiced. This suggests that distraction by changing-state speech and distrac-
 644 tion by whispered speech is the result of two distinct mechanisms: While changing-state
 645 speech is more disruptive due to interference with deliberate serial-order processing (i.e.,
 646 interference-by-process [Jones *et al.*, 1996; Marsh *et al.*, 2009](#)), whispered phonation may
 647 cause additional disruption due to attentional capture or by demanding cognitive resources
 648 to process whispered speech in a comprehensible language (i.e., speech decoding and lexical
 649 access, which may be a more automatic / less conscious process compared to attentional
 650 capture). Interestingly, the magnitude of the changing-state effect was larger than the mag-
 651 nitude of the whispering effect ($\hat{\eta}_p^2$ is about 5.5 times larger for the changing-state effect in

652 Experiment 1), suggesting that interference-by-process causes considerably more disruption
653 than attentional capture by specific features of a whispered voice. It is also possible that
654 attentional capture by whispered speech has been partially reduced through the listeners'
655 cognitive control (depending on their available working memory capacity during the task;
656 e.g., [Kattner, 2021](#); [Sörqvist, 2010](#)) (but see [Körner et al., 2017](#)), whereas the more auto-
657 matic interference-by-process probably cannot be reduced at a cognitive level ([Hughes et al.,](#)
658 [2013](#)). There are several possibilities concerning the cues in whispered speech that may cap-
659 ture attention. Whispered speech is often used to convey secret or personal information
660 and it may thus be considered as potentially more important or self-relevant to a listener.
661 Similarly, whispering may indicate social exclusion from a group and thus trigger an emo-
662 tional response that directs attention to whispered speech. However, it is also possible that
663 it just requires more listening effort and thus attentional control to process and understand
664 whispered speech due to the absence of certain phonetic cues (e.g., the periodic excitation
665 pattern and harmonic structure of modal speech).

666 Experiment 2 was conducted to test whether the impairment of serial short-term memory
667 with whispered speech may be related to attentional capture. Therefore, the same German
668 irrelevant speech materials were presented to participants who did not understand the lan-
669 guage. In a foreign and therefore incomprehensible language, participants are not expected
670 to engage in additional listening effort when processing whispered speech than when pro-
671 cessing voiced speech (e.g., in line with the 'framework for understanding effortful listening'
672 or the 'ease of language understanding' account [Pichora-Fuller et al., 2016](#); [Rönnerberg et al.,](#)
673 [2013](#)). Moreover, whispering in an incomprehensible language may not be a useful cue of

674 enhanced importance or self-relevance of the ‘task-irrelevant’ information. Thus, whispered
675 speech would not be expected to capture more attention than voiced speech when presented
676 in an incomprehensible language. In contrast, interference with serial order processing (due
677 to the changing-state nature of speech) should be unaffected by a change of the language
678 (e.g., [Jones *et al.*, 1990](#)). It was found in Experiment 2 that serial recall was disrupted by
679 the presence of changing-state speech (compared to steady-state speech) – though less than
680 in Experiment 1 – but it was not affected by the phonation type of irrelevant speech. Whis-
681 pered speech even tended to be less disruptive than voiced speech, presumably due to the
682 lower level or spectro-temporal variation reducing interference with serial-order processing
683 (compare [Alikadic and Röer, 2022](#)). The results of Experiment 2 appear to rule out the
684 notion that whispered speech produces greater disruption than normally phonated speech
685 due to the triggering of affective responses - since these should arguably transcend language.

686 The greater disruption produced by whispered against normally phonated speech for
687 native language listeners coheres with previous findings, demonstrating that meaningful
688 sentences produce greater disruption of serial recall than incomprehensible degraded speech
689 or random sequences of spoken syllables, presumably due to higher familiarity or interest
690 (e.g., [Hughes and Marsh, 2020](#); [Kattner *et al.*, 2022](#)). Moreover, the results also gel with
691 the finding that ignoring a telephone conversation whereby only one of the two speakers
692 was audible produced more disruption to a visually-based task than the same conversation
693 wherein both speakers could be heard ([Emberson *et al.*, 2010](#); [Marsh *et al.*, 2018](#)). Similarly,
694 it has been reported that disruption decreases with an increasing number of voices in multi-
695 speaker speech babble background situations ([Jones and Macken, 1995](#); [Zaglauer *et al.*,](#)

696 2017). All three instances (whispered speech, single-sided telephone conversations and multi-
697 speaker background babble) contain intelligible/semi-intelligible speech that involuntarily
698 engages the listener’s attention due to the semantic content (or the potential for meaning).
699 All three types of speech engage cognitive processes related to understanding language, even
700 when the listener is trying to focus on an unrelated, visual task. The same pattern did
701 not emerge when the speech was incomprehensible (Marsh *et al.*, 2018), suggesting that the
702 semantic properties of the half conversation generated a “need to listen” or “involuntary
703 eavesdropping”. It is possible that whispered speech, meaningful to the listener, provokes
704 a similar mechanism of attentional diversion. In contrast to previous findings, however,
705 whispered stimuli do not have to be semantically rich to attract attention. In Experiment
706 1 the whispering effect was similar in magnitude for sequences comprising a repeated single
707 word (e.g., “hand”) as it was for multi-word semantically rich sentences (e.g., prose). Thus,
708 it would appear that lexical identification of a single-item is sufficient to drive the additional
709 disruption produced by whispering as compared to normally phonated speech (possibly due
710 to increased cognitive demand for successful decoding and lexical access, see Pichora-Fuller
711 *et al.*, 2016; Rönnberg *et al.*, 2013).

712 The notion that the whispering effect emerges due to the recruitment of more listening
713 effort for lexical-semantic identification of whispered speech, implies that similar disruptive
714 effects should be observed for speech that is rendered slightly less intelligible via other means
715 of acoustical manipulation. However, such a pattern appears to be absent from previous
716 studies in which the degree of disruption in serial recall typically declines with continuous
717 degradation of the speech signal (lower numbers of frequency bands in noise-vocoded speech

718 or longer segment durations in locally time-reversed speech; see [Ellermeier *et al.*, 2015](#); [Ueda](#)
719 [et al., 2019](#)). This may be because there is an optimal level of intelligibility required for the
720 recruitment of listening effort, or because some other factor (e.g., socio-affective; [Cirillo and](#)
721 [Todt, 2005](#); [Laver, 1994](#)) provokes greater listening effort.

722 Disruption of serial recall in the present experiments was also related to certain psychoa-
723 coustic properties of the irrelevant speech sounds, particularly sound pressure level, loudness,
724 fluctuation strength, and tonality. However, in contrast to previous reports of louder sounds
725 being equally or more disruptive ([Alikadic and Röer, 2022](#); [Ellermeier and Hellbrück, 1998](#)),
726 recall accuracy (not distraction) increased with both the sound pressure level and the loud-
727 ness of the irrelevant speech samples – in particular in Experiment 1. This finding most
728 likely reflects the fact that whispered speech, when presented in a comprehensible language,
729 was more disruptive despite its lower intensity and reduced vowel amplitudes compared to
730 voiced speech (but see [Hughes *et al.*, 2005a](#)). In future work, it may be worth validating
731 whether the disruptive effects of whispered speech in a comprehensible language hold true
732 when whispering is presented at more realistic playback levels. Similarly, irrelevant speech
733 with higher fluctuation strength also led to higher recall accuracy in Experiment 1, but
734 not in Experiment 2. Hence, in line with other previous findings ([Ellermeier *et al.*, 2015](#),
735 also observing higher recall accuracy in the conditions with maximum fluctuation strength),
736 fluctuation strength alone does not seem to be an appropriate predictor of auditory dis-
737 traction (i.e., in the present experiments, it was associated with less distraction; in contrast
738 to [Schlittmeier *et al.*, 2012](#)). A similar relationship was observed also between roughness
739 and memory performance, with voiced speech being characterized by higher roughness – in

740 particular for the male voice – which was associated with higher recall accuracy when the lan-
741 guage was comprehensible (Experiment 1), but with lower accuracy when the language was
742 incomprehensible (Experiment 2). Finally, higher tonality was also associated with higher
743 accuracy in the serial recall task, indicating that the tonality of voiced speech (characterized
744 by more pronounced harmonics) does not necessarily produce more disruption in a serial
745 recall task. Specifically, it appears that certain cues in comprehensible whispered speech
746 (e.g., semantics or social-cognitive aspects) capture attention and produce even more dis-
747 ruption compared to the acoustically driven interference due to changes in vowel amplitudes
748 in voiced speech.

749 While the results of Experiments 1 and 2 support the notion that the changing-state
750 effect and the whispering effect are underpinned by distinct cognitive mechanisms, further
751 studies could add weight to the proposed dichotomy of distraction effects. Previous research
752 suggests that the changing-state effect occurs most prominently in tasks drawing on serial
753 rehearsal (e.g., [Beaman and Jones, 1997](#)), whereas attentional capture effects should arise
754 for any cognitively demanding task (e.g., [Vachon et al., 2017](#)). If the whispering effect for
755 native language listeners is indeed attributable to attentional diversion then it should also
756 be observed on focal tasks that do not draw upon serial processing, such as the missing-item
757 task ([Hughes et al., 2007](#); [Jones and Macken, 1993](#)). Further, the whispering effect unlike
758 the changing-state effect should be influenced by extrinsic or intrinsic cognitive control. For
759 example, the magnitude of the whispering effect for native listeners should be reduced under
760 high task-encoding or cognitive load (see [Hughes et al., 2013](#); [Marsh et al., 2020, 2018](#)) and
761 for individuals with higher working memory capacity (which reflects a trait capacity for

762 cognitive/attentional control; [Hughes *et al.*, 2013](#); [Marsh *et al.*, 2017](#); [Sörqvist *et al.*, 2012](#)).
763 Furthermore the whispering effect should be reduced by previous exposure to distractors
764 (i.e., foreknowledge) which has been shown to reduce the additional disruption produced by
765 comprehensible over incomprehensible spoken sentences ([Kattner *et al.*, 2022](#)) and emotional
766 (e.g., taboo) over neutral words ([Rettie *et al.*, 2024](#)), through reducing the personal relevance,
767 interest ([Hughes and Marsh, 2020](#); [Kattner *et al.*, 2022](#)), or affective responses ([Rettie *et al.*,](#)
768 [2024](#)) produced by the stimuli.

769 V. CONCLUSION

770 Taken together the present study shows that task-irrelevant whispered speech can be
771 more – not less – disruptive to cognitive performance when the language is comprehensible
772 to the listener, but not when the listeners did not understand the language. This suggests
773 that certain semantic and/or social-cognitive features conveyed by whispered voices may
774 capture attention or encourage enhanced listening effort, leading to a lack of cognitive re-
775 sources being available for the focal short-term memory task. In line with an attentional
776 interpretation of the whispering effects, distraction did not increase with psychoacoustic
777 loudness or fluctuation strength – as would have been predicted by a unitary interference-
778 by-process account of auditory distraction (i.e., whispered speech is softer, less tonal, and
779 less fluctuating and should therefore cause less interference with serial-order processing). At
780 the same time, it was observed that changing-state speech is more disruptive than steady-
781 state speech regardless of whether the phonation was voiced or whispered. This indicates
782 two functionally distinct and additive mechanisms of distraction, with one being based on

783 interference between auditory grouping (of changing-state sounds) and deliberate seriation
784 processes, and the other being based on attentional capture by whispered voices.

785 The findings of the current study may have significant practical implications across var-
786 ious real-world contexts, particularly when the whispered speech is intelligible to listeners.
787 While whispering is often employed in open office settings to lower conversational volume
788 for politeness and to avoid disturbing others, as well as to communicate sensitive informa-
789 tion (Cirillo, 2004; Cirillo and Todt, 2005), our study reveals that intelligible whispered
790 speech may attract more attention and lead to increased errors and reduced efficiency com-
791 pared to voiced speech. Over time, these declines in cognitive performance could result
792 in substantial costs for businesses, highlighting the need for sound management strategies
793 or workspace redesigns, such as designated quiet zones, to minimize disruptive intelligible
794 whispers. The results of the present study suggest that whispered conversations can be more
795 distracting than previously thought, calling into question the effectiveness of such policies
796 in maintaining a focused environment. Similarly, in educational environments, intelligible
797 whispers might hinder classroom learning and academic achievement, necessitating strict
798 noise management during study and exam sessions. Whispering during lectures, a common
799 and allegedly unproblematic behavior, may negatively influence academic performance of
800 other students and/or disrupt the lecturer even more than voiced conversations. Also the
801 concept of ‘whisper zones’ in libraries, which are intended to minimize disruption with read-
802 ing, may be based on a misguided assumption. In high-stakes cognitive settings, such as
803 hospital operating rooms or control rooms, whispering may disrupt critical tasks, underscor-
804 ing the importance of stringent sound management policies in these areas. Moreover, while

805 whispered speech is commonly used to maintain privacy in public or semi-public spaces
806 like libraries, our research suggests that intelligible whispers may undermine this intent and
807 disrupt others more than anticipated. Organizations should encourage staff to reconsider
808 the use of whispered conversations and explore physical barriers or soundproofing options to
809 enhance privacy and minimize disruption. Finally, the tendency for intelligible whispering to
810 be more disruptive than voiced speech has implications for individuals with breathy or soft
811 voices, whether due to natural variations or clinical conditions characterized by breathiness
812 (e.g., vocal fold paralysis; [Macdonell and Holmes, 2007](#)). Such individuals may inadvertently
813 create more disruptions than those with clearer vocal tones, especially if their speech is eas-
814 ily understood. This highlights the potential need for voice training or sound management
815 strategies in shared environments.

816 VI. AUTHOR DECLARATIONS

817 A. Conflict of Interest

818 The authors have no conflicts of interest to declare.

819 B. Ethics Approval

820 The two experiments in this study were conducted strictly in accordance with the Ethical
821 Principles of the Acoustical Society of America and the Declaration of Helsinki. All par-
822 ticipants were informed about the duration and procedure, potential risks, data protection
823 regulations, and their right to withdraw from participating at any time, without conse-

824 quence. Written informed consent was obtained prior to the start of the experiment. The
825 experimental protocols were approved by the ethics committee of the University of Lincoln
826 (ref: 33415).

827 **VII. DATA AVAILABILITY**

828 The data and analysis scripts of the experiments in this study are available upon request
829 from the corresponding author.

830 **ACKNOWLEDGMENTS**

831 We are thankful to our student research assistants Leonardo Stuff, Ramona Spitz, Patrizia
832 Scholz, Cosima Stokar von Neuforn, Marko Sanden, and Maike Gerlach for their help with
833 the production of speech recordings, the recruitment of participants and the data collection of
834 Experiment 1. We also thank Lia Downing, Veronika Schaer, Mia Dennis, and Joshua Dudley
835 for their help with the recruitment of participants and the data collection for Experiment 2.

836 **VIII. REFERENCES**

837

838 Alikadic, L., and Röer, J. P. (2022). “Loud Auditory Distractors Are More Difficult to Ignore
839 after All: A Preregistered Replication Study with Unexpected Results,” *Experimental*
840 *Psychology* **69**(3), 163–171, doi: [10.1027/1618-3169/a000554](https://doi.org/10.1027/1618-3169/a000554).

841 Beaman, C. P., and Jones, D. M. (1997). “Role of serial order in the irrelevant speech effect:
842 Tests of the changing-state hypothesis,” *Journal of Experimental Psychology: Learning*
843 *Memory and Cognition* **23**(2), 459–471, doi: [10.1037/0278-7393.23.2.459](https://doi.org/10.1037/0278-7393.23.2.459).

844 Bell, R., Buchner, A., and Mund, I. (2008). “Age-Related Differences in Irrelevant-Speech
845 Effects,” *Psychology and Aging* **23**(2), 377–391, doi: [10.1037/0882-7974.23.2.377](https://doi.org/10.1037/0882-7974.23.2.377).

846 Bell, R., Mieth, L., Röer, J. P., and Buchner, A. (2022). “The metacognition of au-
847 ditory distraction: Judgments about the effects of deviating and changing auditory
848 distractors on cognitive performance,” *Memory & Cognition* 2021 **50**(1), 1–14, doi:
849 [10.3758/s13421-021-01200-2](https://doi.org/10.3758/s13421-021-01200-2).

850 Bell, R., Röer, J. P., Dentale, S., and Buchner, A. (2012). “Habituation of the irrelevant
851 sound effect: Evidence for an attentional theory of short-term memory disruption,” *Journal*
852 *of Experimental Psychology: Learning Memory and Cognition* **38**(6), 1542–1557, doi: [10.](https://doi.org/10.1037/a0028459)
853 [1037/a0028459](https://doi.org/10.1037/a0028459).

854 Bell, R., Röer, J. P., Lang, A. G., and Buchner, A. (2019). “Distraction by steady-
855 state sounds: Evidence for a graded attentional model of auditory distraction,” *Journal*
856 *of Experimental Psychology: Human Perception and Performance* **45**(4), 500–512, doi:

857 [10.1037/xhp0000623](https://doi.org/10.1037/xhp0000623).

858 Benjamini, Y., and Hochberg, Y. (1995). “Controlling the false discovery rate: A practical
859 and powerful approach to multiple testing,” *Journal of the Royal Statistical Society, Series*
860 *B* **57**(1), 289 – 300, doi: [10.2307/2346101](https://doi.org/10.2307/2346101).

861 Bregman, A. S. (1990). *Auditory Scene Analysis: The perceptual organization of sound*
862 (MIT Press, Cambridge, MA).

863 Campbell, T., Beaman, C. P., and Berry, D. C. (2002). “Auditory memory and the irrelevant
864 sound effect: Further evidence for changing-state disruption,” *Memory* **10**(3), 199–214, doi:
865 [10.1080/09658210143000335](https://doi.org/10.1080/09658210143000335).

866 Cirillo, J. (2004). “Communication by unvoiced speech: the role of whispering,” *Annals of*
867 *the Brazilian Academy of Sciences* **76**, 413–423, doi: [10.1590/S0001-37652004000200034](https://doi.org/10.1590/S0001-37652004000200034).

868 Cirillo, J., and Todt, D. (2005). “Perception and judgement of whispered vocalisations,”
869 *Behaviour* **142**(1), 113–129, doi: [10.1163/1568539053627758](https://doi.org/10.1163/1568539053627758).

870 Colle, H. A., and Welsh, A. (1976). “Acoustic masking in primary memory,” *Journal of Ver-*
871 *bal Learning and Verbal Behavior* **15**(1), 17–31, doi: [10.1016/S0022-5371\(76\)90003-7](https://doi.org/10.1016/S0022-5371(76)90003-7).

872 Cowan, N. (1995). *Attention and Memory: An Integrated Framework* (Oxford University
873 Press).

874 DIN Deutsches Institut für Normung e.V. (2009). “Measurement technique for the simula-
875 tion of the auditory sensation of sharpness,” Standard DIN 45692:2009-08.

876 DIN Deutsches Institut für Normung e.V. (2010). “Calculation of loudness level and loud-
877 ness from the sound spectrum - zwicker method - amendment 1: Calculation of the loudness
878 of time-variant sound; with cd-rom,” Standard DIN 45631/A1:2010-03.

- 879 Dorsi, J., Viswanathan, N., Rosenblum, L. D., and Dias, J. W. (2018). “The role
880 of speech fidelity in the irrelevant sound effect: Insights from noise-vocoded speech
881 backgrounds,” *Quarterly Journal of Experimental Psychology* **71**(10), 2152–2161, doi:
882 [10.1177/1747021817739257](https://doi.org/10.1177/1747021817739257).
- 883 Eckert, H., and Laver, J. (1994). *Menschen und ihre Stimmen: Aspekte der vokalen Kom-*
884 *munikation* (Beltz/Psychologie Verlags Union).
- 885 ECMA International (2022). “Psychoacoustic metrics for itt equipment - part 2 (models
886 based on human perception),” Standard ECMA 418-2 (2nd edition).
- 887 Eimer, M., Nattkemper, D., Schröger, E., and Prinz, W. (1996). “Involuntary attention,”
888 in *Handbook of Perception and Action*, edited by O. Neumann and F. Sanders (Academic
889 Press, London, UK), Chap. 5, pp. 389–446, doi: [10.1016/S1874-5822\(96\)80022-3](https://doi.org/10.1016/S1874-5822(96)80022-3).
- 890 Eklund, I., and Traunmüller, H. (1997). “Comparative Study of Male and Female Whispered
891 and Phonated Versions of the Long Vowels of Swedish,” *Phonetica* **54**(1), 1–21, doi: [10.1159/000262207](https://doi.org/10.1159/000262207) publisher: De Gruyter Mouton.
- 893 Ellermeier, W., and Hellbrück, J. (1998). “Is Level Irrelevant in ”Irrelevant Speech”?
894 Effects of Loudness, Signal-to-Noise Ratio, and Binaural Unmasking,” *Journal of Ex-*
895 *perimental Psychology: Human Perception and Performance* **24**(5), 1406–1414, doi:
896 [10.1037/0096-1523.24.5.1406](https://doi.org/10.1037/0096-1523.24.5.1406).
- 897 Ellermeier, W., Kattner, F., Ueda, K., Doumoto, K., and Nakajima, Y. (2015). “Memory
898 disruption by irrelevant noise-vocoded speech: Effects of native language and the number
899 of frequency bands,” *The Journal of the Acoustical Society of America* **138**(3), 1561–1569,
900 doi: [10.1121/1.4928954](https://doi.org/10.1121/1.4928954).

- 901 Ellermeier, W., and Zimmer, K. (1997). “Individual differences in susceptibility to the
 902 ‘irrelevant speech’ effect,” *Journal of the Acoustical Society of America* **102**, 2191–2199,
 903 doi: [10.1121/1.419596](https://doi.org/10.1121/1.419596).
- 904 Ellermeier, W., and Zimmer, K. (2014). “The psychoacoustics of the irrelevant sound effect,”
 905 *Acoustical Science and Technology* **35**, 10–16, doi: [10.1250/ast.35.10](https://doi.org/10.1250/ast.35.10).
- 906 Emberson, L. L., Lupyan, G., Goldstein, M. H., and Spivey, M. J. (2010). “Overheard
 907 Cell-phone Conversations: When Less Speech is More Distracting,” *Psychological Science*
 908 **21**(10), 1383–1388, doi: [10.1177/0956797610382126](https://doi.org/10.1177/0956797610382126).
- 909 Fastl, H. (1982). “Fluctuation strength and temporal masking patterns of amplitude-
 910 modulated broadband noise,” *Hearing Research* **8**(1), 59–69, doi: [10.1016/
 911 0378-5955\(82\)90034-X](https://doi.org/10.1016/0378-5955(82)90034-X).
- 912 Fastl, H., and Zwicker, E. (2007). *Psychoacoustics: Facts and models*, 3rd ed. ed. (Springer,
 913 Heidelberg, Germany).
- 914 Frühholz, S., Trost, W., and Grandjean, D. (2016). “Whispering - The hidden side of
 915 auditory communication,” *NeuroImage* **142**, 602–612, doi: [10.1016/j.neuroimage.2016.
 916 08.023](https://doi.org/10.1016/j.neuroimage.2016.08.023).
- 917 Genuit, K. (1996). “Objective Evaluation of Acoustic Quality Based on a Relative Ap-
 918 proach,” in *Inter-Noise’96, 25th Anniversary Congress Liverpool*, pp. 1061 p1 – 1061 p6.
- 919 Günther, F., Müller, H. J., and Geyer, T. (2017). “Salience, attention, and perception,”
 920 in *Entrenchment and the psychology of language learning: How we reorganize and adapt
 921 linguistic knowledge*, Language and the human lifespan series (De Gruyter Mouton, Boston,
 922 MA, US), pp. 289–312, doi: [10.1037/15969-014](https://doi.org/10.1037/15969-014).

- 923 Hadlington, L., Bridges, A. M., and Darby, R. J. (2004). “Auditory location in the irrelevant
924 sound effect: The effects of presenting auditory stimuli to either the left ear, right ear or
925 both ears,” *Brain and Cognition* **55**, 545–557, doi: [10.1016/j.bandc.2004.04.001](https://doi.org/10.1016/j.bandc.2004.04.001).
- 926 Heeren, W. F. L. (2015). “Vocalic correlates of pitch in whispered versus normal speech,”
927 *The Journal of the Acoustical Society of America* **138**(6), 3800–3810, doi: [10.1121/1.](https://doi.org/10.1121/1.4937762)
928 [4937762](https://doi.org/10.1121/1.4937762).
- 929 Hughes, R. W. (2014). “Auditory distraction: A duplex-mechanism account,” *PsyCh Jour-*
930 *nal* **3**(1), 30–41, doi: [10.1002/pchj.44](https://doi.org/10.1002/pchj.44).
- 931 Hughes, R. W., Hurlstone, M. J., Marsh, J. E., Vachon, F., and Jones, D. M. (2013).
932 “Cognitive control of auditory distraction: impact of task difficulty, foreknowledge, and
933 working memory capacity supports duplex-mechanism account.,” *Journal of experimental*
934 *psychology. Human perception and performance* **39**(2), 539–553, doi: [10.1037/a0029064](https://doi.org/10.1037/a0029064).
- 935 Hughes, R. W., and Marsh, J. E. (2019). “Dissociating two forms of auditory distraction in
936 a novel Stroop serial recall experiment,” *Auditory Perception and Cognition* **2**(3), 129–142.
- 937 Hughes, R. W., and Marsh, J. E. (2020). “When is forewarned forearmed? Predicting au-
938 ditory distraction in short-term memory,” *Journal of Experimental Psychology: Learning*
939 *Memory and Cognition* **46**(3), 427–442, doi: [10.1037/xlm0000736](https://doi.org/10.1037/xlm0000736).
- 940 Hughes, R. W., Tremblay, S., and Jones, D. M. (2005a). “Disruption by speech of serial
941 short-term memory: The role of changing-state vowels,” *Psychonomic Bulletin and Review*
942 **12**, 886–890, doi: [10.3758/BF03196781](https://doi.org/10.3758/BF03196781).
- 943 Hughes, R. W., Vachon, F., and Jones, D. M. (2005b). “Auditory attentional capture
944 during serial recall: Violations at encoding of an algorithm-based neural model?,” *Journal*

- 945 of Experimental Psychology: Learning, Memory, and Cognition **31**(4), 736–749, doi: [10.1037/0278-7393.31.4.736](https://doi.org/10.1037/0278-7393.31.4.736).
- 946
- 947 Hughes, R. W., Vachon, F., and Jones, D. M. (2007). “Disruption of short-term memory by
- 948 changing and deviant sounds: Support for a duplex-mechanism account of auditory dis-
- 949 traction.,” *Journal of Experimental Psychology: Learning, Memory, and Cognition* **33**(6),
- 950 1050–1061, doi: [10.1037/0278-7393.33.6.1050](https://doi.org/10.1037/0278-7393.33.6.1050).
- 951 Ito, T., Takeda, K., and Itakura, F. (2005). “Analysis and recognition of whispered speech,”
- 952 *Speech Communication* **45**(2), 139–152, doi: [10.1016/j.specom.2003.10.005](https://doi.org/10.1016/j.specom.2003.10.005).
- 953 Jones, D. M., Alford, D., Macken, W. J., Banbury, S. P., and Tremblay, S. (2000). “Inter-
- 954 ference from degraded auditory stimuli: Linear effects of changing-state in the irrelevant
- 955 sequence,” *The Journal of the Acoustical Society of America* **108**(3), 1082–1088, doi:
- 956 [10.1121/1.1288412](https://doi.org/10.1121/1.1288412).
- 957 Jones, D. M., Beaman, C. P., and Macken, W. J. (1996). “The object-oriented episodic
- 958 record model,” in *Models of short-term memory*, edited by S. E. Gathercole (Psychology
- 959 Press, Hove), pp. 209–238.
- 960 Jones, D. M., and Macken, W. J. (1993). “Irrelevant tones produce an irrelevant speech
- 961 effect: Implications for phonological coding in working memory,” *Journal of Experimental*
- 962 *Psychology: Learning, Memory, and Cognition* **19**, 369–381, doi: [10.1037/0278-7393.](https://doi.org/10.1037/0278-7393.19.2.369)
- 963 [19.2.369](https://doi.org/10.1037/0278-7393.19.2.369).
- 964 Jones, D. M., and Macken, W. J. (1995). “Auditory Babble and Cognitive Efficiency: Role
- 965 of Number of Voices and Their Location,” *Journal of Experimental Psychology: Applied*
- 966 **1**(3), 216–226, doi: [10.1037/1076-898X.1.3.216](https://doi.org/10.1037/1076-898X.1.3.216).

- 967 Jones, D. M., Macken, W. J., and Murray, A. C. (1993). “Disruption of visual short-
968 term memory by changing-state auditory stimuli: The role of segmentation,” *Memory &*
969 *Cognition* **21**, 318–328, doi: [10.3758/BF03208264](https://doi.org/10.3758/BF03208264).
- 970 Jones, D. M., Madden, C., and Miles, C. (1992). “Privileged access by irrelevant speech to
971 short-term memory: The role of changing state,” *The Quarterly Journal of Experimental*
972 *Psychology* **44A**, 645–669.
- 973 Jones, D. M., Miles, C., and Page, J. (1990). “Disruption of proofreading by irrelevant
974 speech: Effects of attention, arousal or memory?,” *Applied Cognitive Psychology* **4**, 89–
975 108, doi: [10.1002/acp.2350040203](https://doi.org/10.1002/acp.2350040203).
- 976 Jones, D. M., and Tremblay, S. (2000). “Interference in memory by process or content? A
977 reply to Neath (2000),” *Psychonomic Bulletin and Review* **7**(3), 550–558, doi: [10.3758/
978 BF03214370](https://doi.org/10.3758/BF03214370).
- 979 Jovičić, S. T., and Šarić, Z. (2008). “Acoustic Analysis of Consonants in Whispered Speech,”
980 *Journal of Voice* **22**(3), 263–274, doi: [10.1016/j.jvoice.2006.08.012](https://doi.org/10.1016/j.jvoice.2006.08.012).
- 981 Kattner, F. (2021). “Transfer of working memory training to the inhibitory control of audi-
982 tory distraction,” *Psychological Research* **85**, 1–15, doi: [10.1007/s00426-020-01468-0](https://doi.org/10.1007/s00426-020-01468-0).
- 983 Kattner, F. (2024). “False memories through auditory distraction: When irrelevant speech
984 produces memory intrusions in the absence of semantic interference,” *Quarterly Journal*
985 *of Experimental Psychology* doi: [10.1177/17470218241235654](https://doi.org/10.1177/17470218241235654).
- 986 Kattner, F., and Bryce, D. (2022). “Attentional control and metacognitive monitor-
987 ing of the effects of different types of task-irrelevant sound on serial recall,” *Journal*
988 *of Experimental Psychology: Human Perception and Performance* **48**(2), 139–158, doi:

989 [10.1037/xhp0000982](https://doi.org/10.1037/xhp0000982).

990 Kattner, F., and Ellermeier, W. (2014). “Irrelevant speech does not interfere with serial
991 recall in early blind listeners,” *Quarterly Journal of Experimental Psychology* **67**(11),
992 2207–2217, doi: [10.1080/17470218.2014.910537](https://doi.org/10.1080/17470218.2014.910537).

993 Kattner, F., and Ellermeier, W. (2018). “Emotional prosody of task-irrelevant speech in-
994 terferes with the retention of serial order,” *Journal of Experimental Psychology: Human*
995 *Perception and Performance* **44**(8), 1303–1312, doi: [10.1037/xhp0000537](https://doi.org/10.1037/xhp0000537).

996 Kattner, F., Fischer, M., Caling, A. L., Cremona, S., Ihle, A., Hodgson, T., and Föcker, J.
997 (2024). “The disruptive effects of changing-state sound and emotional prosody on verbal
998 short-term memory in blind, visually impaired, and sighted listeners,” *Journal of Cognitive*
999 *Psychology* **36**, 28–41, doi: [10.1080/20445911.2023.2186771](https://doi.org/10.1080/20445911.2023.2186771).

1000 Kattner, F., Hanl, S., Paul, L., and Ellermeier, W. (2023). “Task-specific auditory distrac-
1001 tion in serial recall and mental arithmetic,” *Memory and Cognition* **51**(4), 930–951.

1002 Kattner, F., Richardson, B. H., and Marsh, J. E. (2022). “The Benefit of Foreknowledge
1003 in Auditory Distraction Depends on the Intelligibility of pre-exposed Speech,” *Auditory*
1004 *Perception & Cognition* **5**(3-4), 151–168, doi: [10.1080/25742442.2022.2089525](https://doi.org/10.1080/25742442.2022.2089525).

1005 Konno, H. (2016). “Analysis on Acoustical and Perceptual Characteristics of Whispered
1006 Speech and Whisper-to-Normal Speech Conversion” doi: [10.14943/doctoral.k12482](https://doi.org/10.14943/doctoral.k12482).

1007 Körner, U., Röer, J. P., Buchner, A., and Bell, R. (2017). “Working memory capacity is
1008 equally unrelated to auditory distraction by changing-state and deviant sounds,” *Journal*
1009 *of Memory and Language* **96**, 122–137, doi: [10.1016/j.jml.2017.05.005](https://doi.org/10.1016/j.jml.2017.05.005).

- 1010 Labonté, K., Marsh, J. E., and Vachon, F. (2022). “Distraction by auditory semantic devi-
 1011 ations is unrelated to working memory capacity: Further evidence of a distinction between
 1012 acoustic and categorical deviation effects,” *Auditory Perception & Cognition* .
- 1013 Laver, J. (1994). Cambridge Textbooks in Linguistics *Principles of Phonetics* (Cambridge
 1014 University Press), pp. 190–192.
- 1015 LeCompte, D. C., Neely, C. B., and Wilson, J. R. (1997). “Irrelevant speech and irrel-
 1016 evant tones: The relative importance of speech to the irrelevant speech effect,” *Jour-
 1017 nal of Experimental Psychology: Learning Memory and Cognition* **23**(2), 472–483, doi:
 1018 [10.1037/0278-7393.23.2.472](https://doi.org/10.1037/0278-7393.23.2.472).
- 1019 Ljungberg, J. K., Parmentier, F. B., Hughes, R. W., Macken, W. J., and Jones, D. M.
 1020 (2012). “Listen Out! Behavioural and Subjective Responses to Verbal Warnings,” *Applied
 1021 Cognitive Psychology* **26**(3), 451–461, doi: [10.1002/acp.2818](https://doi.org/10.1002/acp.2818).
- 1022 Macdonell, R. A., and Holmes, R. (2007). *Motor Speech and Swallowing Disorders*, 155–170
 1023 (Elsevier), doi: [10.1016/B978-0-323-03354-1.50016-X](https://doi.org/10.1016/B978-0-323-03354-1.50016-X).
- 1024 Marsh, J. E., Campbell, T., Vachon, F., Taylor, P., and Hughes, R. W. (2020). “How the
 1025 deployment of visual attention modulates auditory distraction,” *Attention, Perception,
 1026 and Psychophysics* **82**, 350–362, doi: [10.3758/s13414-019-01800-w](https://doi.org/10.3758/s13414-019-01800-w).
- 1027 Marsh, J. E., Hughes, R. W., and Jones, D. M. (2009). “Interference by process, not
 1028 content, determines semantic auditory distraction,” *Cognition* **110**(1), 23–38, doi: [10.
 1029 1016/j.cognition.2008.08.003](https://doi.org/10.1016/j.cognition.2008.08.003).
- 1030 Marsh, J. E., Vachon, F., and Sörqvist, P. (2017). “Increased distractibility in schizotypy:
 1031 Independent of individual differences in working memory capacity?,” *Quarterly Journal of*

- 1032 Experimental Psychology **70**(3), 565–578, doi: [10.1080/17470218.2016.1172094](https://doi.org/10.1080/17470218.2016.1172094).
- 1033 Marsh, J. E., Yang, J., Qualter, P., Richardson, C., Perham, N., Vachon, F., and Hughes,
1034 R. W. (2018). “Postcategorical auditory distraction in short-term memory: Insights from
1035 increased task load and task type.,” *Journal of Experimental Psychology: Learning, Mem-
1036 ory, and Cognition* **44**(6), 882–897, doi: [10.1037/xlm0000492](https://doi.org/10.1037/xlm0000492).
- 1037 Miles, C., Jones, D. M., and Madden, C. A. (1991). “Locus of the Irrelevant Speech Effect
1038 in Short-Term Memory,” *Journal of Experimental Psychology: Learning, Memory, and
1039 Cognition* **17**(3), 578–584, doi: [10.1037/0278-7393.17.3.578](https://doi.org/10.1037/0278-7393.17.3.578).
- 1040 Peirce, J. W., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., Kast-
1041 man, E., and Lindeløv, J. K. (2019). “PsychoPy2: Experiments in behavior made easy,”
1042 *Behavior Research Methods* **51**, 195–203, doi: [10.3758/s13428-018-01193-y](https://doi.org/10.3758/s13428-018-01193-y).
- 1043 Pichora-Fuller, M. K., Kramer, S. E., Eckert, M. A., Edwards, B., Hornsby, B. W., Humes,
1044 L. E., Lemke, U., Lunner, T., Matthen, M., Mackersie, C. L., Naylor, G., Phillips, N. A.,
1045 Richter, M., Rudner, M., Sommers, M. S., Tremblay, K. L., and Wingfield, A. (2016).
1046 “Hearing impairment and cognitive energy: The framework for understanding effortful
1047 listening (fuel),” *Ear and Hearing* **37**, 5S–27S, doi: [10.1097/AUD.0000000000000312](https://doi.org/10.1097/AUD.0000000000000312).
- 1048 Rettie, L., Potter, R. F., Brewer, G., Degno, F., Vachon, F., Hughes, R. W., and Marsh, J. E.
1049 (2024). “Warning—taboo words ahead! Avoiding attentional capture by spoken taboo dis-
1050 tractors,” *Journal of Cognitive Psychology* **36**(1), 61–77, doi: [10.1080/20445911.2023.
1051 2285860](https://doi.org/10.1080/20445911.2023.2285860).
- 1052 Röer, J. P., Bell, R., and Buchner, A. (2013). “Self-relevance increases the irrelevant sound
1053 effect: Attentional disruption by one’s own name,” *Journal of Cognitive Psychology* **25**(8),

- 1054 925–931, doi: [10.1080/20445911.2013.828063](https://doi.org/10.1080/20445911.2013.828063).
- 1055 Röer, J. P., Bell, R., and Buchner, A. (2015). “Specific foreknowledge reduces auditory
1056 distraction by irrelevant speech,” *Journal of Experimental Psychology: Human Perception*
1057 *and Performance* **41**(3), 692–702, doi: [10.1037/xhp0000028](https://doi.org/10.1037/xhp0000028).
- 1058 Röer, J. P., Körner, U., Buchner, A., and Bell, R. (2017a). “Attentional capture by taboo
1059 words: A functional view of auditory distraction,” *Emotion* **17**(4), 740–750, doi: [10.1037/
1060 emo0000274](https://doi.org/10.1037/emo0000274).
- 1061 Röer, J. P., Rummel, J., Bell, R., and Buchner, A. (2017b). “Metacognition in Auditory
1062 Distraction: How Expectations about Distractibility Influence the Irrelevant Sound Effect,”
1063 *Journal of Cognition* **1**(1), 2, doi: [10.5334/joc.3](https://doi.org/10.5334/joc.3).
- 1064 Rönningberg, J., Holmer, E., and Rudner, M. (2021). “Cognitive hearing science: three
1065 memory systems, two approaches, and the ease of language understanding model,”
1066 *Journal of Speech, Language, and Hearing Research* **64**, 359–370, doi: [10.1044/2020_
1067 JSLHR-20-00007](https://doi.org/10.1044/2020_JSLHR-20-00007).
- 1068 Rönningberg, J., Lunner, T., Zekveld, A., Sörqvist, P., Danielsson, H., Lyxell, B., Örjan
1069 Dahlström, Signoret, C., Stenfelt, S., Pichora-Fuller, M. K., and Rudner, M. (2013). “The
1070 ease of language understanding (elu) model: Theory, data, and clinical implications,”
1071 *Frontiers in Systems Neuroscience* **7**, 48891, doi: [10.3389/fnsys.2013.00031](https://doi.org/10.3389/fnsys.2013.00031).
- 1072 Salamé, P., and Baddeley, A. D. (1982). “Disruption of short-term memory by unattended
1073 speech: Implications for the structure of working memory,” *Journal of Verbal Learning*
1074 *and Verbal Behavior* **21**, 150–164, doi: [10.1016/S0022-5371\(82\)90521-7](https://doi.org/10.1016/S0022-5371(82)90521-7).

- 1075 Salamé, P., and Baddeley, A. D. (1989). “Effects of background music on phonological short-
 1076 term memory,” *The Quarterly Journal of Experimental Psychology Section A* **41A**(1),
 1077 107–122, doi: [10.1080/14640748908402355](https://doi.org/10.1080/14640748908402355).
- 1078 Schlittmeier, S. J., Hellbrück, J., and Klatte, M. (2008). “Does irrelevant music cause an
 1079 irrelevant sound effect for auditory items?,” *European Journal of Cognitive Psychology*
 1080 **20**, 252–271, doi: [10.1080/09541440701427838](https://doi.org/10.1080/09541440701427838).
- 1081 Schlittmeier, S. J., Weißgerber, T., Kerber, S., Fastl, H., and Hellbrück, J. (2012).
 1082 “Algorithmic modeling of the irrelevant sound effect (ISE) by the hearing sensation
 1083 fluctuation strength,” *Attention, Perception, and Psychophysics* **74**(1), 194–203, doi:
 1084 [10.3758/s13414-011-0230-7](https://doi.org/10.3758/s13414-011-0230-7).
- 1085 Sörqvist, P. (2010). “High working memory capacity attenuates the deviation effect but not
 1086 the changing-state effect: Further support for the duplex-mechanism account of auditory
 1087 distraction,” *Memory and Cognition* **38**(5), 651–658, doi: [10.3758/MC.38.5.651](https://doi.org/10.3758/MC.38.5.651).
- 1088 Sörqvist, P., Nöstl, A., and Halin, N. (2012). “Working memory capacity modulates habit-
 1089 uation rate: Evidence from a cross-modal auditory distraction paradigm,” *Psychonomic*
 1090 *Bulletin and Review* **19**, 245–250, doi: [10.3758/s13423-011-0203-9](https://doi.org/10.3758/s13423-011-0203-9).
- 1091 Tremblay, S., and Jones, D. M. (1999). “Change of intensity fails to produce an irrele-
 1092 vant sound effect: Implications for the representation of unattended sound,” *Journal of*
 1093 *Experimental Psychology: Human Perception and Performance* **25**(4), 1005–1015, doi:
 1094 [10.1037/0096-1523.25.4.1005](https://doi.org/10.1037/0096-1523.25.4.1005).
- 1095 Tremblay, S., Nicholls, A. P., Alford, D., and Jones, D. M. (2000). “The Irrelevant Sound
 1096 Effect: Does Speech Play a Special Role?,” *Journal of Experimental Psychology: Learning*

- 1097 Memory and Cognition **26**(6), 1750–1754, doi: [10.1037/0278-7393.26.6.1750](https://doi.org/10.1037/0278-7393.26.6.1750).
- 1098 Ueda, K., Nakajima, Y., Kattner, F., and Ellermeier, W. (2019). “Irrelevant speech effects
1099 with locally time-reversed speech: Native vs non-native language,” The Journal of the
1100 Acoustical Society of America **145**(6), 3686, doi: [10.1121/1.5112774](https://doi.org/10.1121/1.5112774).
- 1101 Vachon, F., Labonté, K., and Marsh, J. E. (2017). “Attentional capture by deviant sounds:
1102 A noncontingent form of auditory distraction?,” Journal of Experimental Psychology:
1103 Learning Memory and Cognition **43**(4), 622–634, doi: [10.1037/xlm0000330](https://doi.org/10.1037/xlm0000330).
- 1104 Viswanathan, N., Dorsi, J., and George, S. (2014). “The role of speech-specific proper-
1105 ties of the background in the irrelevant sound effect,” Quarterly Journal of Experimental
1106 Psychology **67**(3), 581–589, doi: [10.1080/17470218.2013.821708](https://doi.org/10.1080/17470218.2013.821708).
- 1107 Wingfield, A. (2016). “Evolution of models of working memory and cognitive resources,”
1108 Ear and Hearing **37**, 35S–43S, doi: [10.1097/AUD.0000000000000310](https://doi.org/10.1097/AUD.0000000000000310).
- 1109 Zaglauer, M., Drotleff, H., and Liebl, A. (2017). “Background babble in open-plan offices:
1110 A natural masker of disruptive speech?,” Applied Acoustics **118**, 1–7, doi: [10.1016/J.
1111 APACOUST.2016.11.004](https://doi.org/10.1016/J.APACOUST.2016.11.004).
- 1112 Zeamer, C., and Fox Tree, J. E. (2013). “The process of auditory distraction: Disrupted
1113 attention and impaired recall in a simulated lecture environment,” Journal of Experimental
1114 Psychology: Learning Memory and Cognition **39**(5), 1463–1472, doi: [10.1037/a0032190](https://doi.org/10.1037/a0032190).
- 1115 Zhang, C., and Hansen, J. H. L. (2018). “Advancements in whispered speech detection for
1116 interactive/speech systems,” in *Signal and Acoustic Modeling for Speech and Communica-*
1117 *tion Disorders* (De Gruyter), pp. 9–32, doi: [10.1515/9781501502415-002](https://doi.org/10.1515/9781501502415-002).

1118 Żygis, M., Pape, D., Koenig, L. L., Jaskuła, M., and Jesus, L. M. (2017). “Segmental cues
1119 to intonation of statements and polar questions in whispered, semi-whispered and normal
1120 speech modes,” *Journal of Phonetics* **63**, 53–74, doi: [10.1016/j.wocn.2017.04.001](https://doi.org/10.1016/j.wocn.2017.04.001).