

1 **Effectiveness of a loudness model for time-varying sounds in equating the loudness of**
2 **sentences subjected to different forms of signal processing**

3

4

Tudor-Cătălin Zorilă and Yannis Stylianou

5

Toshiba Research Europe Ltd., Cambridge Research Laboratory,

6

208 Cambridge Science Park, Milton Road, Cambridge CB4 0GZ, UK

7

catalin.zorila@crl.toshiba.co.uk, yannis.stylianou@crl.toshiba.co.uk

8

9

Sheila Flanagan and Brian C.J. Moore^{a)}

10

Department of Experimental Psychology, University of Cambridge, Downing Street,

11

Cambridge CB2 3EB, UK

12

saf31@cam.ac.uk, bcjm@cam.ac.uk

13

14 ^{a)} Author to whom correspondence should be addressed

15

16

16 **Abstract:** A model for the loudness of time-varying sounds [B.R. Glasberg and B.C.J. Moore
17 (2012). *J. Audio. Eng. Soc.* **50**, 331-342] was assessed for its ability to predict the loudness
18 of sentences that were processed to either decrease or increase their dynamic fluctuations. In a
19 paired-comparison task, subjects compared the loudness of unprocessed and processed
20 sentences that had been equalized in: (1) root-mean square (RMS) level; (2) the peak long-
21 term loudness predicted by the model; (3) the mean long-term loudness predicted by the
22 model. Method 2 was most effective in equating the loudness of the original and processed
23 sentences.

24

25 PACS numbers: 43.66Cb, 43.66Ba

26 **1. Introduction**

27 There has been considerable interest in recent years in the development of methods of
28 processing speech so as to enhance its intelligibility when background noise and/or
29 reverberation are added after the processing has been applied (Yoo *et al.*, 2007; Zorila *et al.*,
30 2012; Cooke *et al.*, 2013). Such methods have potential applications in public address
31 systems and in classrooms for use with special populations, such as children with “auditory
32 processing disorder” (Moore *et al.*, 2013). It would be trivial to improve the intelligibility of
33 speech simply by increasing its level, thereby improving the signal-to-noise ratio (SNR).
34 Therefore, processing methods of this type have typically been evaluated under the constraint
35 that the root-mean-square (RMS) level of the speech should be the same before and after
36 processing (Zorila *et al.*, 2012; Cooke *et al.*, 2013). However, what is important in practical
37 applications is that the loudness of the speech should not be increased by the processing; the
38 loudness must be kept within a range that is judged as comfortable by the majority of
39 listeners. Therefore, it may be more appropriate to assess the processing under the constraint
40 that the *loudness* of the speech should be the same before and after processing. Here, we
41 present an evaluation of the accuracy of the loudness model developed by Glasberg and
42 Moore (2002) in equating the loudness of unprocessed and processed speech.

43 Two types of speech processing were used, both of which have been shown to
44 improve the intelligibility of speech when applied prior to the addition of background noise
45 (Cooke *et al.*, 2013). One method decreased the short-term level fluctuations in the speech
46 (Zorila *et al.*, 2012) while the other increased them (Takou *et al.*, 2013; Zorila and Stylianou,
47 2014) relative to those of the original speech. The processed signals were therefore thought to
48 provide a strong test of the accuracy of the loudness model. The model used, called the time-
49 varying-loudness (TVL) model (Glasberg and Moore, 2002), takes a time waveform as its
50 input and generates three forms of time-varying loudness: the instantaneous loudness, which
51 is assumed not to be available for conscious perception; the short-term loudness, which is
52 intended to represent the impression of the loudness of a short segment of the sound, for
53 example a syllable in a sentence; and the long-term loudness (LTL), which is intended to
54 represent the overall loudness of a longer sample of the sound, for example a whole sentence.

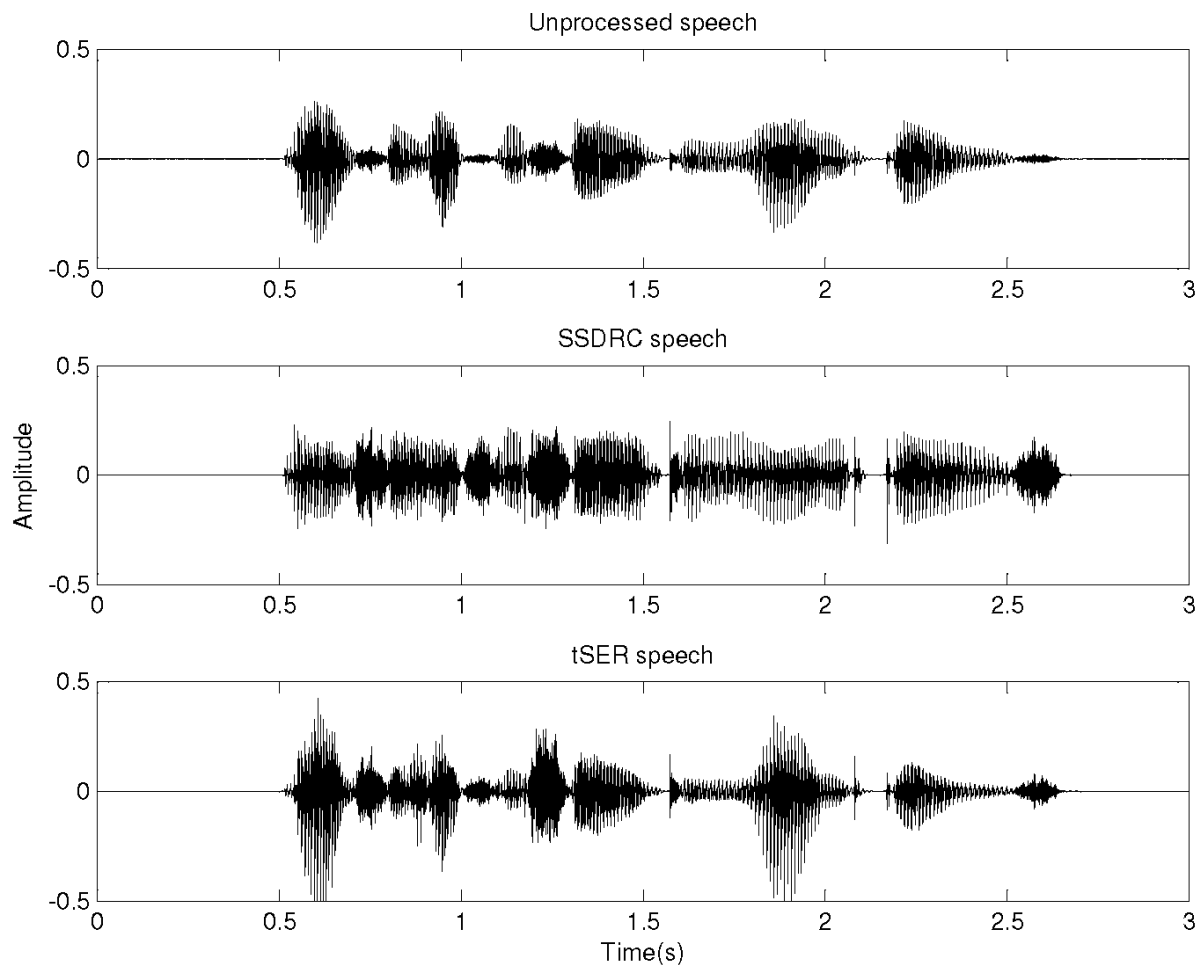
55 In this work, it was also assessed whether overall loudness was best predicted by the
56 maximum value of the LTL reached during presentation of a sentence or by the value of the
57 LTL averaged over all times for which the predicted loudness exceeded a certain threshold
58 value.

59

60 **2. The signal-processing methods**

61 *2.1. Spectral shaping with dynamic range compression*

62 The first method was based on spectral shaping combined with dynamic range compression,
63 denoted SSDRC (Zorila *et al.*, 2012). The signal was analyzed in frames and the spectrum in
64 each frame was estimated by discrete time-frequency transform. Spectral peaks (formants)
65 were sharpened and energy was transferred from low frequencies to medium and high
66 frequencies (1-4 kHz), thereby improving the SNR over the frequency range that is most
67 important for intelligibility (ANSI, 1997). Following spectral shaping, dynamic range
68 compression (DRC) was applied to the broadband signal, aiming to amplify the weaker parts
69 of speech that are more prone to noise masking (fricatives, nasals, and stops), while
70 attenuating parts with more energy (vowels). The effect of the DRC was a reduction of the
71 waveform's envelope variations over time, as illustrated in the middle trace of Fig. 1. Hence
72 the processed speech had a smaller dynamic range than the original speech (top trace).



73

74 Fig. 1. Waveforms of unprocessed speech (top trace), speech processed using SSDRC
75 (middle trace) and speech processed using tSER. All sentences had the same RMS value. The
76 sentence was “Rice is often served in round bowls”.

77

78 2.2. Time-domain spectral energy reallocation

79 The second method was based on reallocation of energy in frequency using time-domain
80 processing, and is denoted tSER (Takou *et al.*, 2013; Zorila and Stylianou, 2014). This had
81 three processing stages. In one stage, the low-frequency components below 400 Hz were
82 isolated by lowpass filtering and were passed on unprocessed for combination with the
83 signals from the other stages. In a second stage, the signal was pre-emphasized with a first-
84 order finite impulse response (FIR) filter that flattened the spectral tilt. The third stage took its
85 input from the second stage and applied a spectral contrast enhancement algorithm
86 resembling the two-tone suppression that occurs in the cochlea (Turicchia and Sarpeshkar,

87 2005). The outputs of all three stages were combined after weighting of their relative
88 magnitudes. The tSER-processed envelope showed increased envelope fluctuations relative to
89 the original speech, and had a greater dynamic range than the original speech, as illustrated in
90 the bottom trace of Fig. 1.

91

92 **3. The loudness model**

93 The TVL model used here (Glasberg and Moore, 2002) was an extension of the model for
94 stationary sounds developed by Moore *et al.* (1997). The transfer of sound through the outer
95 and middle ear was modeled using a single FIR filter. Different filters can be used for
96 different sound presentation methods (e.g., free field, diffuse field, or headphone). Here, the
97 diffuse-field option was used, as the stimuli for the experiment were presented using
98 headphones with a diffuse-field response. The version of the model used here was slightly
99 modified to have the middle-ear transfer function given by Glasberg and Moore (2006), as
100 described by Moore (2014).

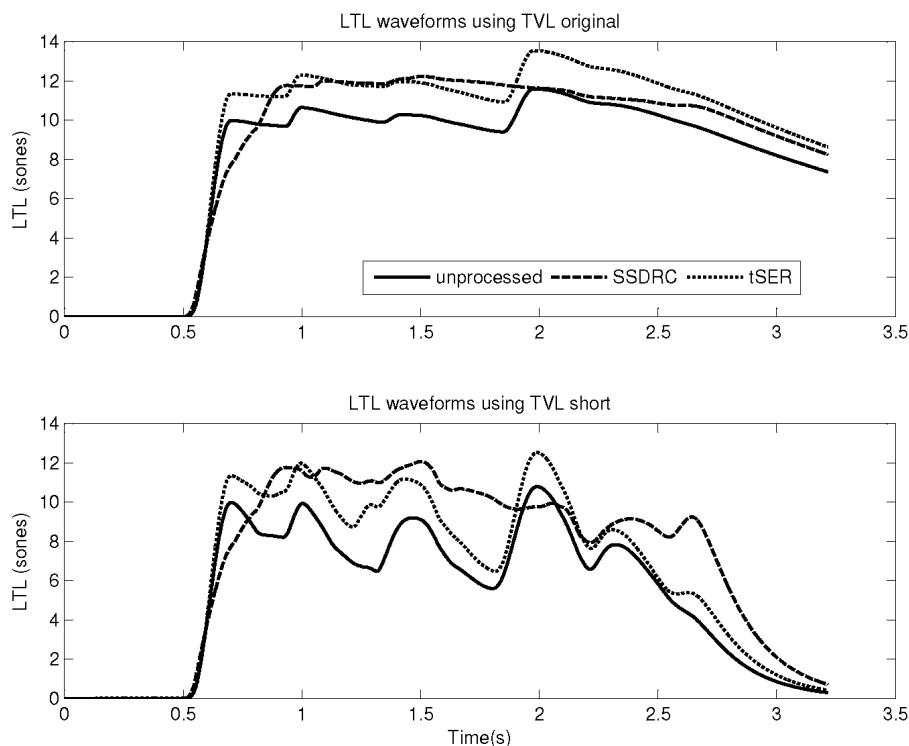
101 A running estimate of the spectrum of the sound at the output of the FIR filter was
102 obtained by calculating six Fast Fourier Transforms (FFTs) in parallel, using signal segment
103 durations that decreased with increasing center frequency. This was done to give sufficient
104 spectral resolution at low frequencies and sufficient temporal resolution at high frequencies.
105 All FFTs were updated every 1 ms. Each FFT was used to calculate spectral magnitudes over
106 a specific frequency range; values outside that range were discarded. An excitation pattern
107 was calculated from the short-term spectrum at 1-ms intervals, using the same method as
108 described by Moore *et al.* (1997). The next stage was the calculation of the “instantaneous”
109 loudness, which is assumed to be an intervening variable that is not available for conscious
110 perception. The calculation of instantaneous loudness from the excitation pattern was done in
111 the same way as described by Moore *et al.* (1997).

112 The short-term loudness was calculated from a running average of the instantaneous
113 loudness, using an averaging process resembling the way that a control signal is generated in
114 an automatic gain control (AGC) circuit. The LTL was calculated from the short-term
115 loudness, again using a form of averaging resembling the operation of an AGC circuit, but

116 with longer time constants. For details, see Glasberg and Moore (2002) and Moore (2014).

117 In the original version of the TVL model, the release time of the averager used to
118 calculate the LTL had a value of 2000 ms. This relatively long time constant was partly
119 chosen to reflect high-level processes such as memory. Here, we evaluated both the original
120 version of the TVL model and a version in which the release time used to calculate the LTL
121 was shorter, at 200 ms.

122 When a single sentence is used as input to the model, the predicted LTL builds up
123 over some time, and then stabilizes at roughly a constant value. However, the value still
124 fluctuates to some extent; the fluctuation is greater when the release time constant is shorter,
125 as illustrated in Fig. 2. The question arises as to whether the overall loudness as judged by
126 human listeners is better predicted by the peak value reached by the LTL or by the mean
127 value of the LTL over the time period where its value is reasonably stable. Both approaches
128 were evaluated here.



129 Fig. 2. Long-term loudness as a function of time predicted by the original TVL model (top)
130 and the version of the model with shorter release time (bottom) for the sentence “Rice is often
131 served in round bowls” either unprocessed (solid lines) or processed using SSDRC (dashed
132 lines) or tSER (dotted lines).

133 **4. Loudness comparison experiments**

134 Two experiments were conducted, one using the original release time to calculate the LTL
135 and one using the shorter release time of 200 ms. These are denoted experiments 1 and 2,
136 respectively.

137

138 *4.1 Subjects*

139 Fifteen subjects (7 male) were tested in experiment 1 and ten subjects (5 male) were tested in
140 experiment 2. All reported having normal hearing and all had audiometric thresholds ≤ 20 dB
141 HL for all audiometric frequencies from 0.25 to 6 kHz. Their ages ranged from 18 to 70 years
142 for both experiments (mean = 40.3 years for experiment 1 and 40.7 years for experiment 2).
143 Five subjects took part in both experiments. All subjects were native speakers of English.

144

145 *4.2 Procedure*

146 A paired-comparison procedure was used. Ten Harvard sentences (Rothausser *et al.*, 1969)
147 were used, spoken by a man. On each trial, the same sentence was presented twice in
148 succession, once unprocessed and once processed with one of the two methods (either
149 SSDRC or tSER). The order of the unprocessed and processed sentences was random with the
150 constraint that the unprocessed sentence occurred equally often in the first and second
151 positions. The unprocessed sentence had an overall diffuse-field equivalent level of 65 dB
152 SPL (its level and spectrum at the eardrum were the same as would be produced if the sound
153 were presented in a diffuse field with a level of 65 dB SPL at the position corresponding to
154 the center of the listener's head). The two sentences within a trial were equalized either in
155 RMS level, in the peak LTL predicted by the TVL model, or in the mean LTL predicted by
156 the TVL model (see below for details of the equalization procedure). The subject was asked
157 to use a slider on a screen, controlled by a computer mouse, to indicate whether the first or
158 the second sentence was louder and by how much. The scale ranged from -3 (sentence 1
159 much louder) to $+3$ (sentence 2 much louder). A slider setting of 0 indicated that the two
160 sentences were equal in loudness. The scale was continuous. All ten sentences were used with
161 each equalization method and processing method. Pairs of sentences for the different

162 equalization methods (3 types) and different speech-processing methods (2 types) were
163 interleaved and presented in an order that was different for each subject, with the constraint
164 that the same sentence was never presented twice in succession.

165 When the unprocessed sentence was judged as louder than the processed sentence, any
166 non-zero response was scored as a negative number. Conversely, when the processed
167 sentence was judged as louder than the unprocessed sentence, any non-zero response was
168 scored as a positive number. The coded responses for each processing method were averaged
169 across all sentences. For simplicity, the result is called the “mean score.” The equalization
170 method that led to a mean score closest to zero was deemed to be the method that gave the
171 most accurate loudness equalization.

172

173 *4.3 Equalization of the original and processed speech*

174 For each unprocessed sentence, the following were calculated: (1) the RMS level; (2) the
175 peak LTL predicted by the TVL model; (3) the mean LTL predicted by the TVL model
176 averaged across all values of the LTL that were above 1 sone. For experiment 1, the overall
177 amplitude of a given processed sentence was iteratively scaled until either: (1) the RMS level
178 was matched to that of the same unprocessed sentence; (2) the peak LTL matched that of the
179 same unprocessed sentence; (3) the mean LTL matched that of the same unprocessed
180 sentence. This was done separately for each sentence. The resulting scaled amplitudes were
181 those used in experiment 1. For experiment 2, the amplitudes of all sentences were scaled
182 either so that the peak LTL was 10 sones (corresponding to the average peak LTL for the
183 unprocessed sentences before scaling) or so that the mean LTL was 7 sones (corresponding to
184 the mean of the mean LTL of the unprocessed sentences before scaling).

185 In experiment 1, for peak LTL equalization, the level of the SSDRC-processed speech
186 was reduced, on average, by 0.2 dB, and that of the tSER-processed speech was reduced by
187 2.6 dB, relative to the levels required for equal RMS. For mean LTL equalization, the level of
188 the SSDRC-processed speech was reduced, on average, by 1.8 dB, while that for tSER-
189 processed speech was reduced by 3.2 dB. Although the mean reduction was very small for
190 peak LTL equalization and SSDRC-processed speech, the change in level varied across

191 sentences from -1.7 to 2.8 dB (standard deviation = 1.4 dB). Thus, equating the RMS level of
192 individual unprocessed and SSDRC-processed sentences would probably not lead to equal
193 loudness for all sentences.

194 In experiment 2, for peak LTL equalization, the level of the SSDRC-processed speech
195 was reduced, on average, by 0.9 dB, and that of the tSER-processed speech was reduced by
196 2.6 dB, relative to the levels required for equal RMS. For mean LTL equalization, the level of
197 the SSDRC-processed speech was reduced, on average, by 4.1 dB, while that for tSER-
198 processed speech was reduced by 2.9 dB.

199 It should be noted that the level reductions described above are not solely a result of
200 differences in the temporal properties of the unprocessed and processed speech; they result at
201 least partly from spectral differences between the unprocessed and processed speech. For a
202 fixed RMS level, both types of processing result in a reduction of low-frequency energy and
203 an increase of medium- and high-frequency energy. The medium and high frequencies
204 contribute more to loudness than the low frequencies, so the spectral changes result in an
205 increase in loudness. This point is discussed in more detail later.

206

207 *4.4 Stimulus generation and presentation*

208 Stimuli were generated digitally (16-bit resolution, 16-kHz sampling rate) and presented via
209 Sennheiser HD580 headphones (Wedemark, Germany), which have approximately a diffuse-
210 field frequency response. Subjects were seated in a sound-attenuating chamber. They
211 responded using a computer mouse, as described above. No feedback was given.

212

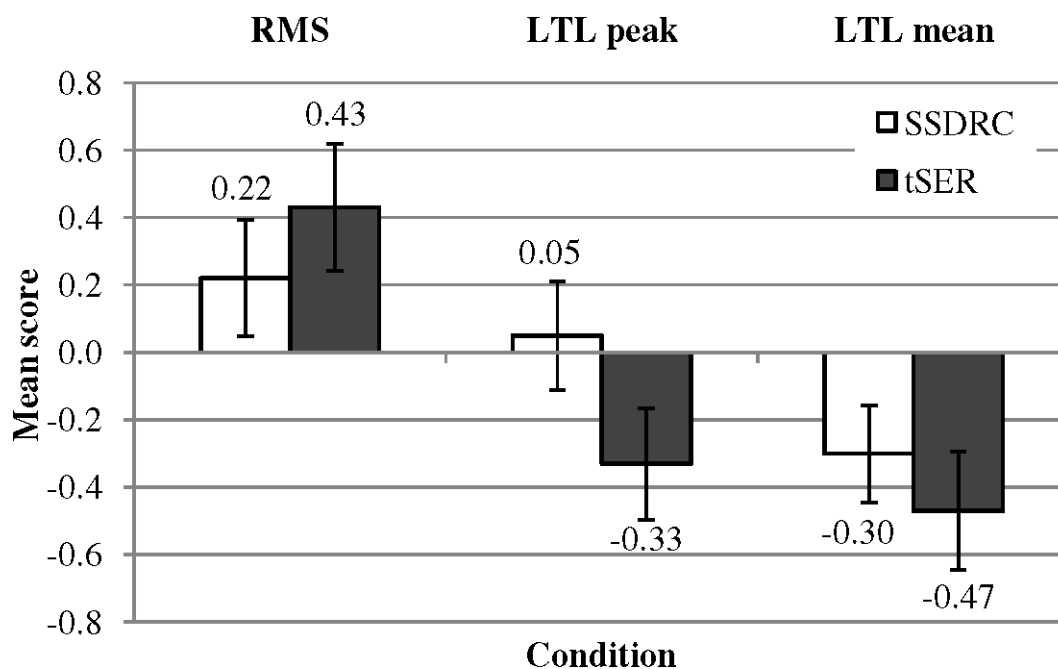
213 **5. Results**

214 *5.1 Experiment 1 (longer release time)*

215 The mean scores for experiment 1, averaged across subjects, are shown in Fig. 3. For
216 sentences equated in RMS level (left pair of bars), the values were positive, by 0.22 scale
217 units for original versus SSDRC and 0.43 scale units for original versus tSER. This means
218 that, at equal RMS level, speech processed using either SSDRC or tSER was louder than the
219 original speech. For sentences equated in peak LTL, the mean score was just above zero

220 (0.05) for original versus SSDRC and below zero (-0.33) for original versus tSER. Thus,
 221 when equated for peak LTL, the processed sentences were well matched in loudness to the
 222 original sentences for SSDRC and were slightly less loud for tSER. For sentences equated in
 223 mean LTL, the mean score was -0.30 for original versus SSDRC and -0.47 for original
 224 versus tSER. Thus, when equated for mean LTL, the tSER-processed sentences were
 225 somewhat less loud than the original sentences. Averaged across processing methods, the
 226 mean scores were 0.32 for RMS equalization, -0.14 for peak LTL equalization, and -0.39 for
 227 mean LTL equalization.

228



229 Fig. 3. Results of experiment 1 (LTL calculated with using the original TVL model with the
 230 longer release time) showing mean ratings of the loudness of processed speech relative to that
 231 of unprocessed speech for two types of processing (SSDRC, open bars, and tSER, shaded
 232 bars) when the unprocessed and processed speech were equated in terms of: (1) RMS level
 233 (left pair of bars); (2) the peak value of the LTL (middle pair of bars); (3) the mean value of
 234 the LTL (right pair of bars). Error bars show ± 1 standard error.

235

236 A two-way repeated-measures analysis of variance (ANOVA) was conducted on the
 237 scores with factors equalization method and type of processing. Mauchly's test indicated that

238 the assumption of sphericity was violated for the factor equalization method and for the
239 interaction of equalisation method with type of processing so the degrees of freedom were
240 adjusted using the Greenhouse-Geisser correction. There was a significant main effect of
241 equalization method: $F(1.07, 14.94) = 23.7, p < 0.001$. There was no significant effect of type
242 of processing ($p > 0.05$), but there was a significant interaction of equalization method and
243 type of processing: $F(1.31, 18.38) = 9.63, p < 0.005$.

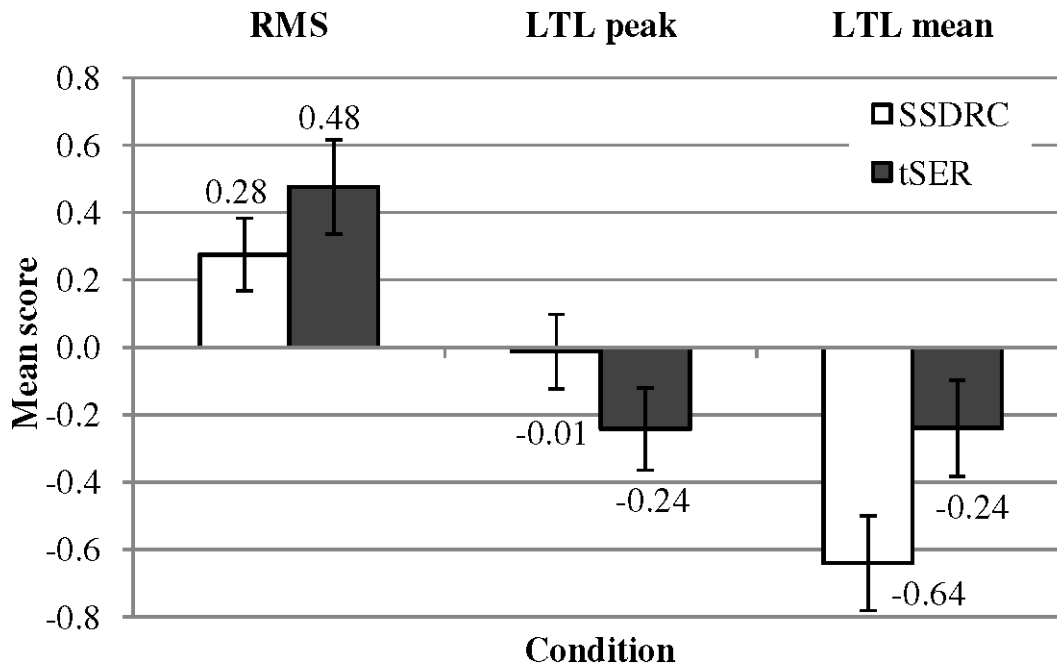
244 A series of *t*-tests (two-tailed) was conducted to assess whether the mean score for
245 each equalization method and processing method was significantly different from zero. For
246 RMS equalization, the mean for tSER was significantly above zero ($t(14) = 2.25, p = 0.041$),
247 but the mean for SSDRC was not. For peak LTL equalization, the means did not differ
248 significantly from zero for either processing method ($p > 0.05$). For mean LTL equalization,
249 the mean for tSER processing was significantly below zero ($t(14) = 2.69, p = 0.018$), while the
250 mean for SSDRC processing did not differ significantly from zero. Overall, it can be
251 concluded that equalization based on the peak LTL was the best method for equating the
252 loudness of the unprocessed and processed sentences.

253

254 5.2 Experiment 2 (shorter release time)

255 The mean scores averaged across subjects are shown in Fig. 4. The pattern of the results is
256 similar to that for experiment 1. For sentences equated in RMS level (left pair of bars), the
257 values were positive, by 0.28 scale units for original versus SSDRC and 0.48 scale units for
258 original versus tSER. Thus, at equal RMS level, speech processed using either SSDRC or
259 tSER was louder than the original speech. For sentences equated in peak LTL, the mean score
260 was very close to zero (-0.01) for original versus SSDRC and slightly below zero (-0.24) for
261 original versus tSER. Thus, when equated for peak LTL, the original and processed sentences
262 were reasonably well matched in loudness. For sentences equated in mean LTL, the mean
263 score was -0.64 for original versus SSDRC and -0.24 for original versus tSER. Thus, when
264 equated for mean LTL, the SSDRC-processed sentences were markedly softer than the
265 original sentences, while the tSER-processed sentences were slightly softer. Averaged across
266 processing methods, the mean scores were 0.38 for RMS equalization, -0.125 for peak LTL

267 equalization, and -0.44 for mean LTL equalization.



268 Fig. 4. As Fig. 3, but showing the results for experiment 2, which used the TVL model with a
 269 shorter release time for calculating the LTL.

270

271 A two-way repeated-measures ANOVA was conducted on the scores with factors
 272 equalization method and type of processing. Mauchly's test indicated that the assumption of
 273 sphericity was violated for the factor equalization method so the degrees of freedom were
 274 adjusted using the Greenhouse-Geisser correction. There was a significant main effect of
 275 equalization method: $F(1.094, 9.825) = 34.8, p < 0.001$. There was no significant effect of
 276 type of processing ($p > 0.05$), but there was a significant interaction of equalization method
 277 and type of processing: $F(2, 18) = 17.5, p < 0.001$.

278 A series of *t*-tests (two-tailed) was conducted to assess whether the mean score for
 279 each equalization method and processing method was significantly different from zero. For
 280 RMS equalization, the means for both SSDRC and tSER were significantly above zero ($t(9) >$
 281 $2.55, p = 0.031$). For peak LTL equalization, the means did not differ significantly from zero
 282 for either processing method ($p > 0.05$). For mean LTL equalization, the mean for SSDRC
 283 processing was significantly below zero ($t(9) = 4.55, p = 0.0014$), while the mean for tSER
 284 processing did not differ significantly from zero. Overall, it can be concluded that

285 equalization based on the peak LTL was the best method for equating the loudness of the
286 unprocessed and processed sentences.

287

288 **6. Discussion**

289 The results showed that for speech processed to either increase or decrease its dynamic range,
290 equalization of the processed and unprocessed speech based on the peak LTL predicted by the
291 TVL model led to more accurate equalization of loudness as perceived by human listeners
292 than equalization based on the mean value of the LTL or the RMS level. This was true for
293 both values of the release time constants used to calculate the LTL.

294 One issue that arises is whether loudness equalization could be performed equally
295 well using a model based on the long-term-average spectra of the stimuli. To assess this, we
296 calculated the long-term spectra across the ten sentences for unprocessed, SSDRC-processed
297 and tSER-processed stimuli with equal RMS levels, and with levels adjusted to give equal
298 peak LTL or equal mean LTL (for the longer time constant only). We then used the loudness
299 model of Moore *et al.* (1997) for stationary sounds, with the modified middle-ear transfer
300 function described by Glasberg and Moore (2006), to predict the loudness in each case. The
301 results are summarized in Table 1. The predicted loudness values for the unprocessed stimuli
302 are all the same, because no adjustment of level was applied for these stimuli. When RMS
303 levels were equalized, the model for stationary sounds predicted that both types of processed
304 stimuli would be louder than the unprocessed stimuli, as found in the data. Equalization based
305 on the mean LTL led to a predicted loudness level that was 1.5 phons higher for SSDRC-
306 processed speech than for unprocessed speech and was almost the same for tSER-processed
307 speech and unprocessed speech. For equalization based on the peak LTL, which led to the
308 most accurate equalization of loudness in our experimental data, the loudness model for
309 stationary sounds did not predict a constant loudness across processing conditions. In
310 particular, the predicted loudness level for SSDRC-processed speech was 2.8 phons higher
311 than for the unprocessed speech and the predicted loudness level for tSET-processed speech
312 was 1 phon higher than for unprocessed speech. We conclude that the dynamic aspects of the
313 stimuli did influence their loudness and that there are benefits in using the model for time-

314 varying sounds to equalize the loudness of unprocessed and processed speech.

315 The results of the experiments are consistent with earlier results showing that the LTL
316 predicted by the TVL model could give accurate predictions of the loudness of speech that
317 had been subjected to multi-channel amplitude compression of the type that is often used in
318 broadcasting (Moore *et al.*, 2003). It is also consistent with the results of Rennie *et al.*
319 (2013), obtained using both loudness matching and categorical loudness scaling, which
320 showed that the LTL gave reasonably accurate predictions of the loudness of a variety of
321 speech-like signals (including speech-shaped noise, unprocessed speech, and speech that was
322 subjected to filtering, reverberation, and amplitude compression and expansion), whereas the
323 predictions were not as accurate when based on the short-term loudness derived using the
324 TVL model or other loudness models (Chalupper and Fastl, 2002; Rennie *et al.*, 2009).

325 Table 2 summarizes the results of the two experiments, showing the mean adjustments
326 in level relative to equal RMS required to equalize the peak LTL or the mean LTL and the
327 mean rating obtained for each equalization method and type of processing. For SSDRC, the
328 correlation between the level adjustments and the ratings for the combined results of the two
329 experiments was 0.97 ($p < 0.05$). The best-fitting linear regression line was

$$330 \quad \text{Rating} = 0.212(\text{Level adjustment}) + 0.18 \quad (1)$$

331 This implies that the rating would be 0, i.e. the SSDRC-processed and unprocessed sentences
332 would be equally loud, when the RMS level of the SSDRC-processed sentences was reduced
333 by 0.8 dB. The level reduction based on equalizing the peak LTL was 0.2 dB with the original
334 release time constant and 0.9 dB with the shorter time constant, both of which are reasonably
335 close to the “ideal” value of 0.8 dB. The level reduction based on equalizing the mean LTL
336 was 1.8 dB with the original release time constant and 4.1 dB with the shorter time constant.
337 This last value is markedly larger than the “ideal” value, suggesting better performance with
338 the original release time.

339 For tSER, the correlation between the level adjustments and the ratings for the
340 combined results of the two experiments was 0.99 ($p < 0.05$). The best-fitting linear
341 regression line was

$$342 \quad \text{Rating} = 0.274(\text{Level adjustment}) + 0.455 \quad (2)$$

343 This implies that the rating would be 0, i.e. the tSER-processed and unprocessed sentences
344 would be equally loud when the RMS level of the tSER-processed sentences was reduced by
345 1.7 dB. The level reduction based on equalizing the peak LTL was 2.6 dB with both the
346 original release time constant and the shorter time constant, reasonably close to the “ideal”
347 value. The level reduction based on equalizing the mean LTL was 3.2 dB with the original
348 release time constant and 2.9 dB with the shorter time constant, both somewhat larger than
349 the “ideal” value.

350 Overall, the results suggest that the level adjustments based on matching the peak
351 value of the LTL were somewhat closer to the adjustments required to actually match the
352 loudness of the unprocessed and processed speech than level adjustments based on the mean
353 value of the LTL. This was the case using both the original release time constant and the
354 shorter time constant. Thus, level adjustments based on matching the peak LTL seem to be
355 preferable.

356

357 **7. Summary and conclusions**

358 Speech processing to enhance its intelligibility when noise is added after processing can
359 either increase or decrease the speech dynamic range, depending on the method of processing,
360 and can also change the average spectral shape of the speech. These changes can alter the
361 loudness of the speech when the overall RMS level is held constant. This paper assessed the
362 effectiveness of three methods in equating the loudness of unprocessed and processed speech,
363 for two methods of speech processing, one that decreased the dynamic range (SSDRC) and
364 one that increased it (tSER). The original and processed speech were equated in terms of: (1)
365 RMS level; (2) the peak LTL predicted by the TVL model; (3) the mean LTL predicted by the
366 TVL model. Two versions of the TVL model were used, one with the original longer release
367 time for calculating the LTL (experiment 1) and the other with a shorter release time
368 (experiment 2).

369 The results were similar for the two experiments. When equated in RMS level, the
370 processed speech was judged as louder than the unprocessed speech for both SSDRC and
371 tSER; the difference was significant for tSER in experiment 1 and for both SSDRC and tSER

372 in experiment 2. When equated in peak LTL, the loudness of the processed speech did not
373 differ significantly from that of the unprocessed speech for either processing method. When
374 equated in mean LTL, the processed speech was judged as softer than the unprocessed
375 speech; the difference was significant for tSER in experiment 1 and for SSDRC in experiment
376 2. It is concluded that the method based on the peak LTL is effective in equating the loudness
377 of processed and unprocessed speech for processing that either decreases or increases the
378 dynamic range of the speech.

379

380 **Acknowledgments**

381 We thank two reviewers for helpful comments on an earlier version of this paper.

382

383 **References**

- 384 ANSI (1997). *ANSI S3.5-1997. Methods for the calculation of the speech intelligibility index*
385 (American National Standards Institute, New York).
- 386 Chalupper, J., and Fastl, H. (2002). "Dynamic loudness model (DLM) for normal and hearing
387 impaired listeners," *Acta Acust. united Ac.* **88**, 378-386.
- 388 Cooke, M., Mayo, C., and Valentini-Botinhao, C. (2013). "Intelligibility-enhancing speech
389 modifications: the Hurricane Challenge," in *Proceedings of Interspeech* (Lyon, France),
390 pp. 3552-3556.
- 391 Glasberg, B. R., and Moore, B. C. J. (2002). "A model of loudness applicable to time-varying
392 sounds," *J. Audio Eng. Soc.* **50**, 331-342.
- 393 Glasberg, B. R., and Moore, B. C. J. (2006). "Prediction of absolute thresholds and equal-
394 loudness contours using a modified loudness model," *J. Acoust. Soc. Am.* **120**, 585-
395 588.
- 396 Moore, B. C. J. (2014). "Development and current status of the "Cambridge" loudness
397 models," *Trends Hear.* **18**, 1-29.
- 398 Moore, B. C. J., Glasberg, B. R., and Baer, T. (1997). "A model for the prediction of
399 thresholds, loudness and partial loudness," *J. Audio Eng. Soc.* **45**, 224-240.
- 400 Moore, B. C. J., Glasberg, B. R., and Stone, M. A. (2003). "Why are commercials so loud? -
401 Perception and modeling of the loudness of amplitude-compressed speech," *J. Audio*
402 *Eng. Soc.* **51**, 1123-1132.

- 403 Moore, D. R., Rosen, S., Bamiou, D. E., Campbell, N. G., and Sirimanna, T. (2013).
404 "Evolving concepts of developmental auditory processing disorder (APD): a British
405 Society of Audiology APD special interest group 'white paper'," *Int. J. Audiol.* **52**, 3-13.
- 406 Rennies, J., Holube, I., and Verhey, J. L. (2013). "Loudness of speech and speech-like
407 signals," *Acta Acust. united Ac.* **99**, 268-282.
- 408 Rennies, J., Verhey, J. L., Chalupper, J., and Fastl, H. (2009). "Modeling temporal effects of
409 spectral loudness summation," *Acta Acust. united Ac.* **95**, 1112-1122.
- 410 Rothausser, E. H., Chapman, W. D., Guttman, N., Nordby, K. S., Silbiger, H. R., Urbanek, G.
411 E., Weinstock, M. (1969). "IEEE recommended practice for speech quality
412 measurements," *IEEE Trans. Audio Electroacoust.* **17**, 225-246.
- 413 Takou, R., Seiyama, N., and Imai, A. (2013). "Improvement of speech intelligibility by
414 reallocation of spectral energy," in *Proceedings of Interspeech* (Lyon, France), pp.
415 3605-3607.
- 416 Turicchia, L., and Sarpeshkar, R. (2005). "A bio-inspired companding strategy for spectral
417 enhancement," *IEEE Trans. Speech. Audio Proc.* **13**, 243-253.
- 418 Yoo, S. D., Boston, J. R., El-Jaroudi, A., Li, C. C., Durrant, J. D., Kovacyk, K., Shaiman, S.
419 (2007). "Speech signal modification to increase intelligibility in noisy environments," *J.*
420 *Acoust. Soc. Am.* **122**, 1138-1149.
- 421 Zorila, C., Kandia, V., and Stylianou, Y. (2012). "Speech-in-noise intelligibility improvement
422 based on spectral shaping and dynamic range compression," in *Proceedings of*
423 *Interspeech* (Portland, OR, USA), pp. 635-638.
- 424 Zorila, C., and Stylianou, Y. (2014). "On spectral and time domain energy reallocation for
425 speech-in-noise intelligibility enhancement," in *Proceedings of Interspeech* (Singapore),
426 pp. 2050-2054.

427

428

428 Table 1. Loudness calculated using a model for stationary sounds based on the long-term
 429 average spectrum, with various forms of equalization across processing method.

430

Equalization method	Unprocessed			SSDRC			tSER		
	RMS	Peak	Mean	RMS	Peak	Mean	RMS	Peak	Mean
Loudness, sones	20.0	20.0	20.0	27.9	24.5	22.3	24.7	21.4	20.7
Loudness level, phons	83.2	83.2	83.2	87.8	86.0	84.7	86.1	84.2	83.7

431

432

432

433

434 Table 2. Summary of the results of experiment 1 (original release time) and experiment 2

435 (shorter release time), showing the average level adjustments (relative to equal RMS)

436 required to equate the peak value of the LTL and the mean value of the LTL, together with

437 the mean loudness ratings.

	RMS		LTL peak		LTL mean	
	SSDRC	tSER	SSDRC	tSER	SSDRC	tSER
Level adjustment re RMS Original, dB	0	0	-0.2	-2.6	-1.8	-3.2
Mean rating	0.22	0.43	0.05	-0.33	-0.30	-0.47
Level adjustment re RMS Shorter, dB	0	0	-0.9	-2.6	-4.1	-2.9
Mean rating	0.28	0.48	-0.01	-0.24	-0.64	-0.24

438