

Western University

Scholarship@Western

Communication Sciences and Disorders
Publications

Communication Sciences and Disorders School

2023

Modified Multiple Stimulus With Hidden Reference and Anchors–Gabrielsson Total Impression Sound Quality Rating Comparisons for Speech in Quiet, Noise, and Reverberation

Mohamed Rahme

Paula Folkeard

Steve Beaulac

Susan Scollie

Vijay Parsa

Follow this and additional works at: <https://ir.lib.uwo.ca/scsdpub>



Part of the [Communication Sciences and Disorders Commons](#)

1 **Modified MUSHRA-Gabrielsson Total Impression Sound**
2 **Quality Rating Comparisons for Speech in Quiet, Noise, and**
3 **Reverberation.**

4 **Mohamed Rahme^{1,2}, Paula Folkeard¹, Steve Beaulac¹, Susan Scollie**
5 **^{1,3} and Vijay Parsa^{1,3}**

6 **Authors' Affiliations:**

7 ¹National Centre for Audiology, Western University, London, Ontario, Canada

8 ²Faculty of Health Sciences, Western University, London, Ontario, Canada

9 ³School of Communication Sciences & Disorders, Western University, London, Ontario, Canada

10 **Corresponding author:**

11 Mohamed Rahme

12 Room 1219, Elborn College Western University London, ON, Canada N6G 1H1. Phone: 519-

13 661-2111. Email: mrahme@uwo.ca

14 **Acknowledgements**

15 The authors would like to thank Muneeb Alam for his assistance with data analysis.

16 **Funding**

17 The Ontario Research Foundation (RE08-072) provided support for this project.

18 **Conflicts of Interest**

19 No conflicts of interest are reported.

20 **Keywords:** MUSHRA, Gabrielsson, sound quality modeling, hearing aids.

21 **ABSTRACT**

22 **Purpose:** The purpose of the study was to obtain, analyze, and compare subjective sound quality
23 data for the same test stimuli using modified multi-stimulus MUSHRA (Multiple Stimulus with
24 Hidden Reference and Anchors) based procedures (*viz.* MUSHRA with custom anchors and
25 MUSHRA without anchor), and the single-stimulus Gabrielsson’s total impression rating
26 procedure. The secondary objective was to evaluate the test-retest reliability of the sound quality
27 procedures

28 **Methods:** Twenty normally hearing young adults were recruited in this study. Participants
29 completed sound quality ratings on two different hearing aid recording datasets – Dataset A
30 contained speech recordings from four different hearing aids under a variety of noisy and
31 processing conditions, and dataset B contained speech recordings from a single hearing aid under
32 a combination of different noisy, reverberant, and signal processing conditions. Recordings in
33 both datasets were rated for their quality using the total impression rating procedure. In addition,
34 quality ratings of dataset A recordings were obtained using a MUSHRA with custom anchors,
35 while the ratings of dataset B recordings were collected using a MUSHRA with no anchors.

36 **Results:** Statistical analyses revealed a high test-retest reliability of quality ratings for the same
37 stimuli that were rated multiple times. In addition, high interrater reliability was observed with
38 all three rating procedures. Further analyses indicated: (i) a high correlation between the total
39 impression rating and the two modified MUSHRA ratings, and (ii) a similar pattern in the
40 variability of the data obtained by the total impression rating and MUSHRA with custom anchors
41 on dataset A, and the total impression rating and the MUSHRA without anchor on dataset B.

42 **Conclusions:** Both sound quality procedures, *viz.* the MUSHRA based procedures and the total
43 impression rating scale, obtained similar quality ratings of varied hearing aid speech recordings
44 with high reliability.

45 **Introduction**

46 Sound quality has been consistently ranked as a top priority for hearing aid users and is found to
47 be the top correlate of consumer satisfaction and preference in both premium and entry-level
48 hearing aids (Kochkin 2010, Saleh et al., 2021; Picou, 2020). To address the quality of sound
49 external to the device, hearing aid manufacturers use digital signal processing strategies such as
50 noise reduction algorithms and directional microphones to help reduce reverberation and
51 background noise. These features improve the signal to noise ratio (SNR), increase speech
52 intelligibility and make listening more comfortable in complex listening environments (Banerjee
53 2011; Cord et al. 2004). However, digital signal processing can generate internal noise, limit
54 available bandwidth, and add distortion to the output signal (Kates & Arehart 2010). These
55 signal degradations can impact sound quality (Arehart et al., 2007 & 2011). To date, hearing aids
56 seem to still underperform in a few listening environments, most prominently in hearing in a
57 classroom, following a conversation in noise, and talking on a phone (Abrams & Kihm, 2015;
58 Picou 2020).

59 Subjective evaluation of sound quality can provide an insight on hearing aid performance
60 as perceived by users for which sound quality continues to play a critical role in hearing aid
61 adoption (Abrams & Kihm, 2015). Therefore, providing optimal sound quality through hearing
62 aid processing is of high interest, and specific methods have been developed for assessing sound
63 quality of hearing aids. The use of magnitude estimation procedures in the context of sound
64 quality and hearing aid literature is well documented (Barry & Kidd 1981; Studebaker &
65 Sherbecoe 1988). An example of such a subjective sound quality procedure is the Gabrielsson
66 scale (Gabrielsson et al., 1988). The Gabrielsson rating scale was initially developed to assess
67 the performance of different types of loudspeakers (Gabrielsson & Sjorgen 1979) and have been

68 used to assess performance of earphones (Gabrielsson et al., 1990), earmolds (Lundberg et al.,
69 1992) and hearing aids (Gabrielsson et al., 1988; Lundberg et al., 1992; Kondo 2012; Narendran
70 & Humes, 2003).

71 During the Gabrielsson rating task, participants are asked to listen to and then rate on a
72 10-point scale the quality of sound signal presented to them. This procedure uses dimension-
73 specific descriptors such as: brightness, clarity, fidelity, fullness, loudness, nearness, softness,
74 spaciousness and total impression. Narendran & Humes (2003) reported that the Gabrielsson task
75 had moderate test-retest reliability across all dimensions for both unaided (median r value= 0.58)
76 and aided (median r value = 0.51) listening conditions. For the purpose of this study, a single
77 dimension “Total Impression” descriptor (taking into consideration the clarity, loudness,
78 background noise, and quality of the recording) was used to evaluate the sound quality of the
79 recordings. This procedure will be referred to as the “total impression rating scale” in this paper.

80 Although using a single rating scale is a fast, direct and subjective evaluation of sound
81 quality, it also has limitations. These caveats may include subjective bias due to context effects,
82 ceiling and floor effects (Studebaker & Sherbecoe 1988). An alternative procedure that is
83 commonly used in the sound quality literature is the Multiple Stimulus with Hidden Reference
84 and Anchors (MUSRHA, ITU- R, 2003). The MUSHRA procedure includes a reference (clean
85 signal), two anchors (distorted signal), and the test stimuli (ITU- R, 2003). Participants begin the
86 task with listening to the reference and then are asked to rate all the other sentences in the group
87 against the reference, in random order. A copy of the reference is included in the randomized list
88 and is referred to as the “hidden” reference. Ratings are provided in a text descriptor along the
89 scale and a percentage score (out of 100%) is then calculated by the software. Ratings of the
90 hidden reference and anchors are included to ensure that the listener uses the full rating scale,

91 and to help prevent ceiling and floor effects. The ratings obtained for the test stimuli may be
92 compared to evaluate various signal processing algorithms.

93 The MUSHRA procedure has been previously applied in assessing the speech and music
94 sound quality of a wide range of signal processing types including, but not limited to, frequency
95 lowering, bandwidth, and noise reduction (Franz & Bitzer, 2010; Glista et al., 2019; Huber et al.,
96 2018; Scollie et al., 2016; Vaisberg et al., 2021). Additionally, MUSHRA's application has been
97 extended to evaluate the sound quality of both speech and music in cochlear implant devices
98 (Caldwell, et al., 2017; Roy et al., 2012).

99 As mentioned previously, the reference signals in standard MUSHRA procedure are
100 clean stimuli with full bandwidth that may be recorded either in quiet or in a positive SNR, while
101 the anchors are traditionally low-pass filtered versions of the original recordings (ITU- R, 2003).
102 However, for assessing the hearing aid sound quality, alternative poor quality or distorted anchor
103 stimuli, such as those containing distortion (e.g., peak or center clipping), reverberation, or
104 background noise, may be more appropriate. For example, Simonsen & Legarth (2010) utilized a
105 recording from an older hearing aid as anchor, in benchmarking the sound quality newer
106 generation hearing aids. Muralimanohar et al. (2013) and Scollie et al. (2016) used noisy speech
107 recordings at low SNRs (-5 dB and -10 dB respectively). This modified MUSHRA based
108 procedure will be referred to for the remainder of this study as MUSHRA with custom anchors.

109 Including references and anchors within the MUSHRA task is part of the standardized
110 procedure (ITU-R, 2003). However, other modification of the standard MUSHRA could be
111 omitting the references and/or anchors. This alternative MUSHRA based procedure may be
112 chosen when sensible reference or anchors cannot be feasibly recorded, or if the underlying
113 constructs associated with "high" or "low" sound quality are not well-understood. For example,

114 Salehi et al. (2018) compared the impact of signal processing features in two different hearing
115 aids on perceived sound quality by listeners with hearing differences by using a modified
116 MUSHRA based procedure that did not have a reference or anchors. This modified MUSHRA
117 procedure will be referred as MUSHRA with no anchor for the remainder of this study.

118 Although MUSHRA has been commonly administered as a subjective sound quality
119 procedure, it requires a computerized interface (Parsa et al., 2013), can be an effortful, relatively
120 costly and time-consuming procedure (Völker et al., 2016). Völker and colleagues argued that
121 the barriers that impact the duration of completing ratings on a traditional MUSHRA dataset can
122 include the factors of age, technological commitment, and hearing status. Furthermore, the
123 sound quality procedures administered in this study, MUSHRA with custom anchors, MUSHRA
124 without anchors, and the total impression rating scale (Gabrielsson et al., 1988) differ in
125 administration and set-up where one method could be advantageous over the other depending on
126 the research question (Studebaker & Sherbecoe 1988; & Völker et al., 2016). The MUSHRA
127 based procedures have the listener compare multiple stimuli on a trial and provide one sound
128 quality rating per stimulus while during the Gabrielsson rating scale procedure, the listener
129 evaluates the sound quality of a single stimulus on a trial and provide ratings along multiple
130 dimensions of sound quality. The rationale for choosing the single descriptor “total impression
131 rating scale” over the rest of the Gabrielsson’s dimensions was to provide a single sound quality
132 rating that is comparable to the MUSHRA based procedures. Narendran & Humes (2003) have
133 shown previously that the single descriptor total impression rating scale was sufficient to judge
134 the sound quality of speech processed in noise and in quiet in an unaided condition.

135 In summary, hearing aid sound quality varies with hearing aid settings and environments
136 and is an important outcome measure that is used in assessing new signal processors. Two sound

137 quality procedures have been used in recent studies: comparisons of multiple stimuli (e.g.,
138 MUSHRA based procedures) and single-stimulus ratings (e.g., total impression rating scale).
139 How these two common procedures relate to one another is not known, it is therefore difficult to
140 compare ratings made with the two procedures. Therefore, the current study is devised to
141 provide initial comparisons between the modified MUSHRA based procedures and the single
142 stimulus rating procedure in a sample of young adults with hearing thresholds within normal
143 range. In particular, the purpose of this study was to evaluate the reliability of, and agreement
144 between multi-stimulus MUSHRA based procedures and the single-stimulus rating procedure for
145 the same test stimuli with a normal hearing young adult participant cohort. The test stimuli were
146 selected across a range of common factors used in hearing aid studies, including signal
147 processing, signal to noise ratio, and reverberation.

148 **Materials and methods**

149 Participants

150 Twenty participants (seven males and thirteen females) were recruited and completed the study
151 protocol. Participant ages ranged from 20-31 years with a mean age of 24.6 (SD 2.95 years).
152 Male ages ranged from 20-24 years (mean 22.7, SD 1.38 years) and female ages ranged from 23-
153 31 (mean 25.72, SD 3.08 years). Participants were native speakers of English, and all
154 participants were tested to confirm that they had audiometric thresholds at 25 dB HL or better
155 between 250 and 8000 Hz, and normal tympanograms prior to starting the study. The
156 experimental procedures were approved by the Western University Human Research Ethics
157 Board (REB 111665). Participants were compensated for their time.

158 Test Stimuli

159 The recordings for this study were made at a constant sound level and frequency content to avoid
160 any potential impact on the dimension-specific ratings. Two sets of stimuli were used in this
161 study (Table 1). These two sets were selected as exemplars of two types of sound quality studies:
162 (1) studies that compare sound quality of a signal processor type *across brands*; and (2) studies
163 that compare variations in signal processing *within one hearing aid*. Both stimulus sets were
164 rated using both a modified MUSHRA procedure and single-stimulus ratings, and both were
165 presented in sound field from zero degrees azimuth at 60 dBA, which was a comfortable level
166 for our normal hearing individuals. Participants were not instructed to change the volume of the
167 presentation level. Both stimulus sets used a range of listening conditions, varying the signal to
168 noise ratio, background noise type, and/or reverberation characteristics of the room to generate
169 stimuli that varied in sound quality. This type of environmental manipulation may be done in
170 hearing aid studies when it is of interest to assess signal processing in a range of acoustic
171 settings. Details of the stimuli are provided for each stimulus set below. The stimuli used in this
172 study came from two previously completed hearing aid projects.

173 MUSHRA with custom anchors stimulus dataset (dataset A)

174 Stimulus set for the MUSHRA with custom anchors procedure included recordings from four
175 different hearing aid brands, generated from a previous study of noise reduction signal
176 processing (Scollie et al., 2016). Briefly, each hearing aid was programmed to a flat moderate
177 hearing loss using DSLv5 child targets. Sentence pairs (IEEE sentences) were recorded through
178 each hearing aid, in three types of background noise [quiet, speech shaped noise (SSN), and
179 babble] at two different SNRs, [0 dB and 5 dB]. Each hearing aid was recorded twice, with the
180 brand-specific digital noise reduction (DNR) signal processing enabled and disabled. For this
181 study, we chose anchors that were the same sentence pairs, mixed with noise at -10 dB SNR

182 using both types of noise, and a reference from one hearing aid recorded in quiet with digital
183 noise reduction off. Two additional stimuli were also included by mixing speech and noise at a
184 +10 dB SNR, to better represent low noise listening conditions in the stimulus set.

185 MUSHRA without anchors stimulus dataset (dataset B)

186 Stimulus set for MUSHRA without anchor procedure used recordings from one, 20 channel
187 digital behind-the-ear style hearing aid. This aid used a hybrid noise management system that
188 combined multiple microphone options, digital noise reduction, and level dependent speech
189 enhancement to create a set of overall noise management profiles ranging from off, through mild,
190 medium, and strong settings of noise management (Table 2). The hearing aid was programmed to
191 a flat moderate hearing loss to DSL v5.0 adult targets and verified using real ear measures. The
192 aids were then placed on a B&K Head and Torso Simulator (HATS) to record hearing aid output.
193 Hearing In Noise Test (HINT) sentences were presented from a desktop computer through a
194 Tucker Davis system and presented through Anthony Gallo Acoustics Nucleus speakers. The
195 HATS was placed in the participant's position in low- and high-reverberation environments. The
196 low-reverberation room was a large double walled sound booth (Industrial Acoustic Company)
197 with a measured reverberation time (RT60) of 100ms. The high reverberation room was a
198 purpose-built reverberation chamber with the dimensions 6.1m x 4.0m x 2.6m. One corner floor
199 to ceiling curtain (135 cm wide) made of micro suede and quilted micro suede and filled with
200 polyester batting and one corner foam pad (122 x 10 x 60cm) were placed in opposite corners of
201 the room. These baffles reduced the reverberation time of the room to 900 ms. The sentence and
202 noise were presented through a desktop computer with levels controlled by a Sound Web and
203 presented through Tannoy Di5 DC speakers.

204 Recordings from these two rooms resulted in stimuli from four digital signal processing
205 settings, two SNR levels (0 dB and 5 dB) and at a several noise conditions [quiet, noise only
206 (SSN and babble), reverberation only, SSN + reverberation and babble + reverberation. The
207 reference was recorded in a quiet environment with minimum sound processing. No additional
208 anchors were included, because the low-SNR high-reverberant conditions were sufficient to
209 represent low sound quality in this stimulus set.

210 Sound Quality Ratings

211 The MUSHRA based procedures were completed on one day and the total impression rating
212 procedure was completed on a separate day. The time period between visits varied
213 between one to seven days with an average of 2.42 days between visits (SD 1.77 days).
214 The starting rating task (i.e., MUSHRA based or total impression rating procedure) was
215 counterbalanced between participants. The startup selection of MUSHRA with custom anchors
216 (dataset A) or MUSHRA without anchors (dataset B) was randomized using an online tool
217 (random.org). The presentation order of the sound files within each rating experiment
218 was then randomized within the custom software that executed the multiple comparison
219 procedure and total impression rating procedure. The rating instructions were provided to each
220 participant in writing and are presented in the supplemental material. Participants completed the
221 experimental procedure on a computer monitor while seated in a double walled sound booth.

222 Ratings using MUSHRA with custom anchors for dataset A

223 The recordings were presented in sound field from zero degrees azimuth. Custom software was
224 developed to mediate the rating data collection. The test stimuli were divided into five sets, with
225 each set representing a noisy condition (quiet, SSN at 0 dB and 5 dB, and babble at 0 dB and 5 dB
226 SNRs). Within each set, there were eleven test stimuli for comparison: the reference recording,

227 recordings from the four hearing aids with the DNR feature on and off, the +10 dB SNR recording,
228 and the anchor (-10 dB SNR recording). The software presented the five sets twice, to gather test-
229 retest ratings for subsequent reliability analyses. In addition, the software randomized the order
230 of the sets, and the order the test stimuli within each set. Figure 1 represents the user interface of
231 the custom MUSHRA software that mediated the data collection.

232 The participants listened to the reference, hidden anchors and test stimuli in each set as
233 many times as they needed to make their ratings, before moving on to the next set. Text descriptors
234 were presented on the software on the left hand side (ranging from “Excellent” to “Bad”). A
235 percentage score appeared on top of the slider as participants moved the slider. The software then
236 logged the percentage scores generated during the rating procedures.

237 Ratings using MUSHRA with no anchors for dataset B

238 The MUSHRA without anchors procedure was administered using the same custom MUSHRA
239 software in a similar manner, where participants were presented with the reference, test stimuli
240 but no anchors (Table 1). For dataset B, there were ten stimulus sets: 5 noise conditions (quiet,
241 SSN 0 dB SNR, SSN 5 dB SNR, Babble 0 dB SNR, Babble 5 dB SNR) x 2 reverberation
242 conditions. Within each set, there was a reference stimulus and the hearing aid recordings under
243 the four signal processing settings (Table 1). On each set, the participants were asked to listen to
244 the reference and then rate the other recordings in comparison to it. Similar to the MUSHRA
245 with custom anchors procedure, the software presented the ten stimulus sets twice, with the order
246 of the sets and order of the test stimuli within the set randomized.

247 Both Stimulus Sets: Total Impression Ratings

248 In this task, the participants were presented with one recording per screen and were asked to
249 move a software slider (Figure 2) to make their rating of the Total Impression of the recording.

250 Participants were asked to take into consideration the clarity, loudness, background noise, and
251 quality of the recording. The participants listened to the recordings as many times as they needed
252 to before rating the total impression of them. Ratings ranged from 0 to 10 with labels on the
253 screen that provided guidance for the raters (e.g., low score indicating poorer sound quality
254 rating and high score indicating better sound quality rating).

255

256 **Results**

257 Reliability of the subjective ratings was analyzed first. The recordings from both the MUSHRA
258 with custom anchors and the MUSHRA without anchor (e.g., references and anchors) that were
259 administered multiple times were examined for their test-retest reliability across all participants.
260 The two-sample Kolmogorov-Smirnov (K-S) test was administered on each participant's test-
261 retest data to ascertain whether they came from the same distribution. For both datasets, the K-S
262 test results showed for 18 of the 20 participants, the two sample K-S test did not reveal that the
263 ratings collected during test and retest sessions were from distinct distributions. To further probe
264 the test-retest reliability, the non-parametric equivalent of the Intraclass Correlation Coefficient
265 (ICC) was employed (Rothery, 1979). The non-parametric statistic was used as the results from
266 the Shapiro-Wilk's test on the test-retest data revealed non-normality. The non-parametric
267 concordance coefficient was computed in R version 4.2.2 using the 'nopaco' package, for all the
268 20 participants and separately for the two datasets. The concordance coefficients ranged
269 between 0.76 to 0.9 for dataset A and between 0.75 to 0.9 for dataset B, with all coefficients
270 statistically significant ($p < 0.025$). As such, the test-retest ratings were averaged for all
271 participants for both MUSHRA based procedures. It is worthwhile to point out that comparable
272 test-retest concordance coefficients were observed for the two MUSHRA based procedures.

273 The histograms of the test-retest averaged rating data in MUSHRA based procedures and the raw
274 Total Impression rating data from all participants across all test conditions in both datasets are
275 depicted in the panels (a) and (b) of Figures 3 and 4 respectively. Deviation from normality is
276 evident in the rating distributions, especially for the Total Impression procedure. The deviation
277 from normality was statistically confirmed through the Shapiro-Wilk's test of normality. The
278 non-parametric concordance coefficient was therefore once again utilized for assessing the inter-
279 rater reliability. The concordance coefficients across MUSHRA with custom anchors and
280 MUSHRA without anchors datasets were 0.69 and 0.70 for the Total Impression ratings, and
281 0.70 and 0.73 for the MUSHRA based procedure ratings respectively, all of which were
282 statistically significant ($p < 0.001$).

283 Figures 3(c) and 4(c) depict the scatter plots between the Total Impression and MUSHRA
284 based procedures both datasets. In order to evaluate how the ratings mapped from one task to the
285 other, a regression analysis was completed using SPSS (SPSS v24; IBM Corporation Chicago,
286 IL). Regression analysis of MUSHRA with custom anchors and MUSHRA without anchors
287 scatter plots revealed correlation coefficients of $r = +0.82$ and $r = +0.85$ respectively. Figures 3(d)
288 and 4(d) display the total impression and MUSHRA based ratings averaged across listeners for
289 each test stimulus in both datasets. The correlation coefficients improved to 0.99 and 0.98
290 respectively, when applied to the condition-averaged data.

291 It is evident from panels (c) and (d) in Figures 3 and 4, there is considerable variability in
292 both the MUSHRA based procedures and the Total Impression ratings, and evidence suggests
293 that this variability must be considered in comparing different rating methods (Höbfeld et al.,
294 2011, Parizet et al., 2005). We followed the procedure outlined for Höbfeld et al., (2011),
295 wherein a quadratic function is fit between the averaged ratings for each test stimulus and the

296 square of their corresponding standard deviations. The quadratic equation is described as: $y^2 =$
297 $\alpha(-x^2 + (\vartheta_{min} + \vartheta_{max})x - \vartheta_{min}\vartheta_{max})$, where x represents the averaged rating, y is the
298 corresponding standard deviation, α is the fitting parameter, and ϑ_{min} and ϑ_{max} are the
299 minimum and maximum rating values respectively. If the rating procedures exhibit similar
300 distribution of ratings across individual participants, then it is expected that they would lead to
301 comparable fitted α parameters (Hoßfeld et al., 2011). We therefore fit separate quadratic
302 functions, parameterized by the α parameter, between the averaged Total Impression or
303 MUSHRA based ratings and their standard deviations, separately for both databases. To derive
304 comparable α parameters, the MUSHRA based ratings were divided by 10, so that both Total
305 Impression and MUSHRA based ratings span between $\vartheta_{min} = 0$ and $\vartheta_{max} = 10$. The α
306 parameter was determined through the least squares fitting function in MATLAB. Figures 5(a)
307 and 5(b) depict the scatter plots of the averaged ratings versus the corresponding standard
308 deviations, along with the fitted quadratic curves, for different rating procedures and across the
309 two datasets. The resulting α parameters were: 0.161 (95% CI: 0.109 – 0.213) and 0.142 (95%
310 CI: 0.1 – 0.185) for the Total Impression Ratings and the MUSHRA with custom anchors
311 ratings; and 0.11 (95% CI: 0.099 – 0.120) and 0.1 (95% CI: 0.092 – 0.108) for the Total
312 Impression ratings and the MUSHRA without anchor ratings. The comparable magnitude of the
313 α parameter highlights the similarity between the speech quality ratings obtained by the Total
314 Impression and MUSHRA based procedures.

315

316 **Discussion**

317 Hearing aid sound quality is an important factor that relates to overall hearing aid outcome and
318 satisfaction (Picou 2020), preference for one hearing aid over another (Saleh et al., 2021) and can
319 be a reason for hearing aid rejection or dissatisfaction if poor (Abrams & Kihm, 2015).

320 Therefore, assessment of hearing aid sound quality is sometimes included in laboratory
321 evaluations of new signal processing (Glista et al., 2019; Huber et al., 2018; Scollie et al., 2016).
322 There are, however, multiple ways to assess sound quality including standardized and multi-
323 stimulus MUSHRA based (ITU- R, 2003) and single-stimulus ratings that may include multiple
324 dimensions of sound quality, based on Gabrielsson's work (Gabrielsson et al., 1988).
325 This study sought to compare the data arising from these two procedures, using a within-subjects
326 design to assess whether and how sound quality ratings differ based on the method used. The
327 findings of this study indicated a high, positive correlation between the ratings of the two
328 procedures. Also, the total impression ratings were positively correlated with both of the multiple
329 comparison datasets.

330 Findings of this study indicated the same stimuli (e.g., multiple references and anchors)
331 that were rated multiple times (during the MUSHRA based procedures) by each participant had a
332 high test-retest reliability meaning that the multiple ratings of the same stimuli were ranked
333 similarly and consistently by participants. Additionally, the multi-stimulus MUSHRA based
334 procedures and the single stimulus rating procedure used in this study showed a high degree of
335 correspondence and agreement in rating the same stimuli. The outcomes of this study suggest
336 that a single-stimulus rating scale, without references and anchors, may be able to generate
337 reliable, consistent and comparable sound quality data when compared to a standardized and
338 commonly used multiple comparison procedures for young, normal-hearing participants.
339 However, the degree to which this result will generalize to other signal processing conditions and
340 other participant populations is unknown and would require further evaluation for conditions not
341 evaluated in this study.

342 Recall that the participants listened to the recordings as many times as they needed to before
343 rating them on the MUSHRA based procedure. This mixed evaluation procedure allowed the
344 participants to compare different recordings in comparison to the reference before finalizing their
345 choice of rating. This methodology is backed up in the literature where a previous study showed
346 a high accuracy in evaluating the pleasantness of different sounds (Parizet et al., 2005).

347 A Standard deviation of Opinion Scores (SOS) analysis was as the Mean Opinion Score, which
348 is the average of the participants' ratings, does not represent rating diversity and therefore only
349 captures a part of the raters' perception. A combined SOS-MOS analysis (Hoßfeld et al., 2011)
350 revealed similar patterns between the total impression ratings and the MUSHRA with custom
351 anchors dataset, as well as the total impression ratings and the MUSHRA without anchor ratings
352 dataset.

353 The MUSHRA without anchor ratings had a higher correlation with the total impression
354 ratings compared to the MUSHRA with custom anchors dataset. Recall that the MUSHRA
355 without anchor dataset used in this study did not have any anchors, which are exemplars of a
356 low-quality sentence pairs. This modification to the MUSHRA procedure was used in this study
357 to extend the comparison to additional stimulus conditions and multiple variations in MUSHRA
358 implementations. The high correlation can be best explained by the fact that the total impression
359 rating scale does not use anchors, just like MUSHRA without anchor dataset, and thus the sound
360 quality ratings of the participants were not influenced by the anchor recordings.

361 The MUSHRA process includes multiple stages of recording the sentence pairs and processing
362 the anchors and test stimuli, such as by peak clipping or asymmetrically distorting the signal.
363 Therefore, setting up the MUSHRA test procedure (or modified versions of it) might require
364 additional efforts, cost and time (Völker et al., 2016). Additionally, participants might be

365 subjected to fatigue due to the extended duration of the test having them listening and rating
366 multiple references, anchors and stimuli. Unlike the single stimulus rating scale, the ratings could
367 be administered on a paper copy (rather than on a computerized interface to accommodate wider
368 range of lab set-ups) or digitally and can be a relatively easy and straightforward procedure.

369 Data collected from sound quality measures are used by hearing aid manufactures and
370 researchers to reinforce and strengthen the objective sound quality models (Völker et al., 2016).
371 Having sound quality data of both subjective and objective procedures in agreement and highly
372 correlated, assists in improving the output of digital signal processing of the HAs. For instance, a
373 previous study indicated a high correlation between the MUSHRA ratings and an objective
374 model of sound quality (Muralimanohar et al., 2013). Findings of the current study demonstrated
375 a high correlation between total impression and MUSHRA based procedures. However, results
376 from MUSHRA without anchors and total impression scales are to be interpreted with caution in
377 the context of reinforcing objective sound quality paradigms. The use of anchors is
378 recommended to reduce any subjective biases that could impact the quality of prediction of
379 mathematical and objective sound quality models (Zielinski et al., 2008).

380 A limitation of this study was recruiting a small sample (N= 20) where all participants
381 were young adults who had normal hearing. As a result, the generalizability of these results to
382 those who are hearing aid users is unknown. Future work is encouraged to consider recruiting
383 participants of different ages with and without hearing loss, and to consider the role of hearing
384 aid experience. The datasets used in this study varied in noise reduction properties, so future
385 studies should assess the sound quality of other signal processing strategies and to investigate
386 how the MUSHRA ratings relate to other dimensions of the Gabrielsson's single-stimulus rating

387 scale. Another limitation of this study was including recordings made at a constant sound level.
388 Future studies are recommended with loudness equalization for the recorded sentence pairs.

389

390 **Conclusion**

391 In conclusion, the current study indicated a high correlation between sound quality data collected
392 via MUSHRA with custom anchors, MUSHRA without anchor and total impression rating
393 scales. Participants reported consistent ratings on the Gabrielsson's total impression scale
394 compared to MUSHRA with custom anchors and MUSHRA without anchor procedures without
395 having to compare the test stimuli to hidden references and anchors. The efficiency and
396 simplicity of administering single stimulus rating scale may facilitate the collection of subjective
397 sound quality data across larger number of participants and wider range of digital signal
398 processing strategies.

399

400 **Data Availability Statement**

401 The datasets generated during and/or analyzed during the current study are not publicly available
402 due to Western University Human Research Ethics Board restriction but are available from the
403 corresponding author on reasonable request.

404

405 **References**

406 Abrams, H. B., & Kihm, J. (2015). An introduction to MarkeTrak IX: A new baseline for the
407 hearing aid market. *The Hearing Review*, 22(6), 16.

408 Arehart K.H., Kates J.M., Anderson M.C. & Harvey L.O. (2007). Effects of noise and distortion
409 on speech quality judgments in normal-hearing and hearing-impaired listeners. *Journal of*
410 *the Acoustical Society of America*, 122, 1150 – 1164.

411 Arehart K.H., Kates J.M. & Anderson M.C. (2010). Effects of noise, nonlinear processing, and
412 linear filtering on perceived speech quality. *Ear & Hearing*, 31, 420 – 436.

413 Arehart, K. H., Kates, J. M., & Anderson, M. C. (2011). Effects of noise, nonlinear processing,
414 and linear filtering on perceived music quality. *International Journal of Audiology*, 50(3),
415 177–190. <https://doi.org/10.3109/14992027.2010.539273>

416 Banerjee, S. (2011). Hearing Aids in the Real World: Typical Automatic Behavior of Expansion,
417 Directionality, and Noise Management, 22, 34–48.

418 Barry, S. J., & Kidd Jr, G. (1981). Psychophysical scaling of distorted speech. *Journal of Speech,*
419 *Language, and Hearing Research*, 24(1), 44-47.

420 Caldwell, M. T., Jiam, N. T., & Limb, C. J. (2017). Assessment and improvement of sound
421 quality in cochlear implant users. *Laryngoscope Investigative Otolaryngology*, 2(3), 119–
422 124. <https://doi.org/10.1002/lio2.71>

423 Cord, M. T., Surr, R. K., Walden, B. E., & Dyrland, O. (2004). Relationship between laboratory
424 measures of directional advantage and everyday success with directional microphone
425 hearing aids. *Journal of the American Academy of Audiology*, 15(05), 353-364.

426 Falk, T. H., Parsa, V., Santos, J. F., Arehart, K., Hazrati, O., Huber, R., Kates, J. M., & Scollie,
427 S. (2015). Objective Quality and Intelligibility Prediction for Users of Assistive Listening
428 Devices: Advantages and limitations of existing tools. *IEEE Signal Processing Magazine*,
429 32(2), 114–124. <https://doi.org/10.1109/MSP.2014.2358871>

430 Franz, S., & Bitzer, J. (2010). Multi-channel algorithms for wind noise reduction and signal
431 compensation in binaural hearing aids. In *Proc. Intl. Workshop Acoust. Echo Noise*
432 *Control (IWAENC)*.

433 Gabrielsson, A., Schenkman, B. N., & Hagerman, B. (1988). The effects of frequency responses
434 on sound quality judgments and speech intelligibility. *Journal of Speech and Hearing*
435 *Research, 31*, 166–177.

436 Gabrielsson, A., Hagerman, B., Bech-Kristensen, T., & Lundberg, G. (1990). Perceived sound
437 quality of reproductions with different frequency responses and sound levels. *The Journal*
438 *of the Acoustical Society of America, 88*(3), 1359–1366. <https://doi.org/10.1121/1.399713>

439 Gabrielsson, A., & Sjögren, H. (1979). Perceived sound quality of sound-reproducing systems.
440 *The Journal of the Acoustical Society of America, 65*(4), 1019–1033.
441 <https://doi.org/10.1121/1.382579>

442 Glista, D., Hawkins, M., Vaisberg, J. M., Pourmand, N., Parsa, V., & Scollie, S. (2019). Sound
443 Quality Effects of an Adaptive Nonlinear Frequency Compression Processor with
444 Normal-Hearing and Hearing-Impaired Listeners. *Journal of the American Academy of*
445 *Audiology, 30*(07), 552–563. <https://doi.org/10.3766/jaaa.16179>

446 Hoßfeld, T., Schatz, R., & Egger, S. (2011). SOS: The MOS is not enough!. In *2011 third*
447 *international workshop on quality of multimedia experience* (pp. 131-136). IEEE.

448 Huber, R., Bisitz, T., Gerkmann, T., Kiessling, J., Meister, H., & Kollmeier, B. (2018).
449 Comparison of single-microphone noise reduction schemes: can hearing impaired
450 listeners tell the difference?. *International Journal of Audiology, 57*(sup3), S55-S61.

451 Huber, R., & Kollmeier, B. (2006). PEMO-Q—A New Method for Objective Audio Quality
452 Assessment Using a Model of Auditory Perception. *IEEE Transactions on Audio, Speech,*

453 *and Language Processing, 14(6), 1902–1911.*

454 <https://doi.org/10.1109/TASL.2006.883259>

455 Huber, R., Parsa, V., & Scollie, S. (2014). Predicting the Perceived Sound Quality of Frequency-

456 Compressed Speech. *PLoS ONE, 9(11), e110260.*

457 <https://doi.org/10.1371/journal.pone.0110260>

458 IEEE. (1969). *IEEE Recommended Practice for Speech Quality Measurements* (Institute of

459 Electrical and Electronic Engineers, New York).

460 ITU-R. (2003). *Recommendation BS.1534: Method for the subjective assessment of*

461 *intermediate quality levels of coding systems.* Geneva, Switzerland: International

462 Telecommunications Union.

463 Kates J.M. & Arehart K.H. (2010). The hearing-aid speech quality index (HASQI). *Journal of*

464 *the Audio Engineering Society, 58, 363 – 381.*

465 Kochkin S. 2010. MarkeTrak VIII: Customer satisfaction with hearing aids is slowly increasing.

466 *Hearing Journal, 63, 11 – 19.*

467 Kondo, K. (2012). *Subjective quality measurement of speech: its evaluation, estimation and*

468 *applications.* Springer Science & Business Media.

469 Lundberg, G. Ovegård, A. Hagerman, B., Gabrielsson, A., & Brändström, U. (1992) Perceived

470 Sound Quality in a Hearing Aid with Vented and Closed Earmould Equalized in

471 Frequency Response, *Scandinavian Audiology, 21(2), 87-*

472 *92, DOI: [10.3109/01050399209045987](https://doi.org/10.3109/01050399209045987)*

473 Moore, B. C. J., Baer, T., Ives, D. T., Marriage, J., & Salorio- Corbetto, M. (2016). Effects of

474 modified hearing aid fittings on loudness and tone quality for different acoustic scenes.

475 *Ear and Hearing, 37(4), 483–491.* <https://doi.org/10.1097/AUD.0000000000000285>

476 Möller, S., Chan, W.-Y., Côté, N., Falk, T. H., Raake, A., & Wältermann, M. (2011). Speech
477 Quality Estimation: Models and Trends. *IEEE Signal Processing Magazine*, 28(6), 18–
478 28. <https://doi.org/10.1109/MSP.2011.942469>

479 Muralimanohar, R. K., Kronen, C., Arehart, K., Kates, J., & Pichora-Fuller, M. K. (2013).
480 Quality of voices processed by hearing aids: Intra-talker differences. *Proceedings of*
481 *Meetings on Acoustics*, 19(1), 060112. <https://doi.org/10.1121/1.4800397>

482 Narendran, M. M., & Humes, L. E. (2003). Reliability and Validity of Judgments of Sound
483 Quality in Elderly Hearing Aid Wearers. *Ear and Hearing*, 24(1), 4–11.
484 <https://doi.org/10.1097/01.AUD.0000051745.69182.14>

485 Nilsson, M., Soli, S. D., & Sullivan, J. A. (1994). Development of the Hearing In Noise Test for
486 the measurement of speech reception thresholds in quiet and in noise. *Journal of the*
487 *Acoustical Society of America*. <https://doi.org/10.1121/1.408469>

488 Parizet, E., Hamzaoui, N., & Sabatie, G. (2005). Comparison of some listening test methods: a
489 case study. *Acta Acustica united with Acustica*, 91(2), 356-364.

490 Parsa, V., Scollie, S., Glista, D., & Seelisch, A. (2013). Nonlinear Frequency Compression:
491 Effects on Sound Quality Ratings of Speech and Music. *Trends in Amplification*, 17(1),
492 54–68. <https://doi.org/10.1177/1084713813480856>

493 Roy, A. T., Jiradejvong, P., Carver, C., & Limb, C. J. (2012). Musical Sound Quality
494 Impairments in Cochlear Implant (CI) Users as a Function of Limited High-Frequency
495 Perception. *Trends in Amplification*, 16(4), 191–200.
496 <https://doi.org/10.1177/1084713812465493>

497 Saleh, H. K., Folkeard, P., Van Eeckhoutte, M., & Scollie, S. (2021). Premium versus entry-level
498 hearing aids: Using group concept mapping to investigate the drivers of preference.

499 *International Journal of Audiology*, 1–15.
500 <https://doi.org/10.1080/14992027.2021.2009923>

501 Salehi, H., Suelzle, D., Folkeard, P., & Parsa, V. (2018). Learning-based reference-free speech
502 quality measures for hearing aid applications. *IEEE/ACM Transactions on Audio,*
503 *Speech, and Language Processing*, 26(12), 2277-2288.

504 Scollie, S., Levy, C., Pourmand, N., Abbasalipour, P., Bagatto, M., Richert, F., Moodie, S.,
505 Crukley, J., & Parsa, V. (2016). Fitting Noise Management Signal Processing Applying
506 the American Academy of Audiology Pediatric Amplification Guideline: Verification
507 Protocols. *Journal of the American Academy of Audiology*, 27(03), 237–251.
508 <https://doi.org/10.3766/jaaa.15060>

509 Studebaker, G. A., & Sherbecoe, R. L. (1988). Magnitude estimations of the intelligibility and
510 quality of speech in noise. *Ear and Hearing*, 9(5), 259-267.

511 Simonsen, C. S., & Legarth, S. V. (2010). A procedure for sound quality evaluation of hearing
512 aids. *Hearing Review*, 17(13), 32-37.

513 Vaisberg, J., Folkeard, P, Levy, S, Dundas, D, Agrawal, S, & Scollie, S. (2021). Sound Quality
514 Ratings of Amplified Speech and Music Using a Direct Drive Hearing Aid: Effects of
515 Bandwidth, *Otology & Neurotology* 42(2), 227-234 doi:
516 10.1097/MAO.0000000000002915

517 Völker, C., Bisitz, T., Huber, R., Kollmeier, B., & Ernst, S. M. A. (2018). Modifications of the
518 MUlti stimulus test with Hidden Reference and Anchor (MUSHRA) for use in audiology.
519 *International Journal of Audiology*, 57(sup3), S92–S104.
520 <https://doi.org/10.1080/14992027.2016.1220680>

521 Zielinski, S., Rumsey, F., & Bech, S. (2008). On some biases encountered in modern audio
522 quality listening tests-a review. *Journal of the Audio Engineering Society*, 56(6), 427-
523 451.

524 **Figure caption**

525 *Figure 1- MUSHRA with custom anchors and MUSHRA without anchor user interface*
526 *screenshot.*

527 *Figure 2. Gabrielsson's Total Impression user interface.*

528 *Figure 3. Comparison of sound quality rating methods for MUSHRA with custom anchors (labeled*
529 *Dataset A). (a) Histogram of the raw ratings across all conditions for the MUSHRA method with*
530 *custom anchors. (b) Histogram of the raw ratings across all conditions from the Total Impression*
531 *method. (c) Scatter plot between the Total Impression and MUSHRA with costumed Anchors*
532 *ratings (MUSHRA-CA). (d) Condition-averaged Total Impression and MUSHRA with costumed*
533 *anchors ratings, where the data was sorted such that the Total Impression ratings ranged from the*
534 *lowest to the highest. Error bars represent the standard deviation of the corresponding mean*
535 *rating.*

536 *Figure 4. Comparison of sound quality rating methods for MUSHRA without anchors (labeled*
537 *Dataset B). (a) Histogram of the raw ratings across all conditions for the MUSHRA without*
538 *anchor method, where no anchors were present. (b) Histogram of the raw ratings across all*
539 *conditions from the Total Impression method. (c) Scatter plot between the Total Impression and*
540 *MUSHRA without anchor ratings (MUSHRA-WA). (d) Condition-averaged Total Impression and*
541 *MUSHRA without anchor ratings for Dataset B, where the data was sorted such that the Total*
542 *Impression ratings ranged from the lowest to the highest. Error bars represent the standard*
543 *deviation of the corresponding mean rating.*

544 *Figure 5 – Scatter plots of the mean sound quality ratings versus their corresponding standard*
545 *deviations, for different rating procedures and datasets. The least-squares curve fits to the scatter*
546 *data are also shown. (a) Dataset A, and (b) Dataset B.*

547
548 **Table caption**

549 *Table 1: Summary of stimulus conditions for the two datasets.*

550 *Table 2: MUSHRA without anchor dataset - Hearing aid settings.*