1

2

3      **Rasch analysis of the Listening Effort Questionnaire – Cochlear Implant (LEQ-CI)**

4      Sarah E. Hughes[1], Alan Watkins[1], Frances Rapport[2], Isabelle Boisvert[3],

5      Catherine M. McMahon[4], Hayley A. Hutchings[1]

6

7

8      [1]Patient and Population Health and Informatics, Swansea University Medical School,

9      Singleton Park, Swansea, UK SA2 8PP

10     [2]Australian Institute of Health Innovation, Macquarie University, NSW, Australia

11     [3]Faculty of Medicine and Health, School of Health Sciences, University of Sydney,

12     Sydney, Australia

13     [4]Audiology, H:EAR, Department of Linguistics, Faculty of Medicine, Health and

14     Human Sciences, Macquarie University, NSW, Australia

15

20

27     **All correspondence should be addressed to:**
28     Sarah Hughes
29     Swansea University Medical School
30     Institute of Life Sciences 2
31     Singleton Park, Swansea, UK, SA2 8PP
32     E-mail: sarah.hughes@swansea.ac.uk
33

34                                                          **ABSTRACT**

35    **Objectives:** Listening effort may be defined as the attentional and cognitive resources needed

36    to understand an auditory message, modulated by motivation. Despite the use of hearing

37    devices such as hearing aids (HAs) or cochlear implants (CIs), the requirement for high

38    listening effort remains a challenge for individuals with hearing loss. The Listening Effort

39    Questionnaire - Cochlear Implant (LEQ-CI) is a hearing-specific patient-reported outcome

40    measure (PROM) that has been designed for use in the CI candidacy and rehabilitation

41    process to assess perceived listening effort in everyday life in adults with severe-profound

42    hearing loss. The LEQ-CI has been developed in line with international consensus-based

43    standards for best practice in PROM construction. The aim of this study was to improve the

44    measurement precision of the LEQ-CI and to assess its psychometric measurement

45    properties.

46    **Design:** A field test was undertaken with 330 CI patients from five NHS auditory implant

47    centres in the United Kingdom (UK). Participants were adults ($\geq$18 years of age), had a

48    severe-profound hearing loss, and met the UK candidacy criteria for cochlear implantation

49    specified by the National Institute for Health and Care Excellence (NICE). Participants

50    completed and returned an anonymised 29-item (each with a 5-point or 7-point response

51    option), draft version of the LEQ-CI (LEQ-CI[29]) and a demographic questionnaire. Rasch

52    analysis was undertaken using Winsteps software and the partial credit model to assess rating

53    scale function and item fit. Results informed refinements to produce a 21-item version (LEQ-

54    CI[21]) which underwent a further Rasch analysis.

55    **Results:** The sample was predominantly female: 60.3% (n = 191). Median age of participants

56    was 66 (range 21 - 89) years, with 7.3% (n = 24) of respondents being CI candidates and

57    92.7% (n = 306) being CI recipients. Mean duration of implantation was 3.8 (SD 4.8) years.

58    Initial Rasch analysis of the LEQ-CI[29] revealed poor rating scale functioning. Collapsing the

59    5- and 7-point rating scales to 3-point and 4-point scales and removing eight items produced

60    a 21-item PROM (LEQ-CI[21]) which met the Rasch criteria for rating scale functioning. Rasch

61    analysis of the LEQ-CI[21] showed good fit to the Rasch model. No items showed misfit and

62    dimensionality analysis supported the existence of a single Rasch dimension, defined as

63    perceived listening effort in daily life. Person reliability was 0.91 and the person separation

64    index was 3.28, establishing four levels of person ability. The item separation index was 9.69,

65    confirming the item hierarchy. No items showed differential item functioning (DIF) for

66    gender or age. The item difficulty range was -0.81 to 1.05, the person ability range for non-

67    extreme persons was -3.54 to 2.49, and the mean person ability was -0.31.

68    **Conclusions:** Overall, the LEQ-CI[21] was found to meet the Rasch model criteria for interval

69    level measurement. The LEQ-CI[21] is the first PROM to be developed specifically for the

70    measurement of perceived listening effort and one of the first PROMs for use with CI

71    patients to be developed using Rasch analysis. The LEQ-CI[21] has the potential to be used as a

72    research tool and in clinical practice to evaluate perceived listening effort in daily life.

73    Further psychometric evaluation of the LEQ-CI[21] is planned. (513words)

74

75                                            **INTRODUCTION**

76          Listening effort may be defined as the mental exertion required to attend to, and

77   understand, an auditory message (McGarrigle et al. 2014). It is a significant problem for

78   individuals with hearing loss insofar as they are required to exert greater listening effort than

79   individuals with normal hearing despite the use of hearing devices such as hearing aids

80   and/or cochlear implants (Ohlenforst & Zekveld 2017). Damage to the auditory pathway

81   associated with sensorineural hearing loss results in the degradation of the quality of the

82   auditory signal as well as its intensity. Furthermore, listening often takes place in less than

83   optimal acoustic conditions requiring individuals to deploy additional cognitive resources to

84   extract, decode and comprehend the incoming auditory message (Edwards 2016). Sustained

85   effortful listening in adults with hearing loss is known to contribute to fatigue (Holman et al.

86   2019; Hornsby 2013; Hornsby et al. 2016; Hornsby & Kipp 2016), increased levels of

87   workplace absence/sickness (Nachtegaal et al. 2009), poor mental health (Amieva et al.

88   2018), social withdrawal and isolation (Ramage-Morin 2016), and negatively impact quality

89   of life (McRackan et al. 2017).

90          A hearing-specific, validated measure of listening effort could provide hearing

91   healthcare professionals with a means to support interventions designed to reduce listening

92   effort and mitigate its wider effects. Several methods have been developed to measure

93   listening effort; however, deployment of these measures has been predominantly within

94   research settings. These techniques include behavioural measures that rely on task

95   performance as a proxy-measure of effort (cf. Gagné et al. 2017 for a review); physiological

96   measures such as electroencephalography (EEG, Miles et al. 2017) and pupillometry

97   (Zekveld et al. 2018); and self-report measures of perceived listening effort. Frequently used

98   alongside behavioural or physiological measures, self-report measures have utilised visual

99    analogue scales (VAS, Bräcker et al., 2019) or ordinal rating scales (Johnson et al. 2015) to

100   rate intensity of effort during a specific listening task.

101            To complicate matters, the relationship between these various measures of listening

102   effort is not yet well understood and there is a generalised lack of agreement in the empirical

103   literature among findings using these different types of measures (Hornsby 2013; Miles et al.

104   2017). This lack of consensus has led to the suggestion that these various methods may be

105   measuring different aspects of the listening effort construct (Strand et al. 2018; Strand et al.

106   2020). For example, Alhanbali et al. (2018) used factor analysis to explore the relationship

107   among behavioural, physiological and self-report measures, revealing four underlying

108   dimensions. It seems likely that a range of measures, with method selection depending on the

109   aims of measurement, will be required if listening effort is to be evaluated appropriately.

110            Patient-reported outcome measures (PROMs) offer promise as a viable clinical tool

111   for clinicians and researchers wishing to gain insight into perceived listening effort in daily

112   life as experienced by individuals with hearing loss. PROMs are self-report measures,

113   typically validated questionnaires, that measure aspects of a person's health (such as

114   symptoms, functioning, and quality of life) known only to the individual (Devlin & Appleby

115   2010; FDA 2020). Originally developed for clinical trials research, these tools are now

116   embedded in routine clinical practice across a range of health conditions (Field et al. 2019).

117   There is a long tradition of self-report in audiology and, commensurate with this tradition,

118   there are a number of hearing-specific PROMs (Granberg et al. 2014). Most of these

119   measures have been developed for use in the population of adults with mild-moderate hearing

120   loss (MMHL) who use hearing aids. Furthermore, most hearing-specific PROMs have been

121   developed and validated using traditional psychometric methods which limits their use as

122   clinical measures with individual patients. Notably, no hearing-specific PROMs have been

123   identified that measure perceived listening effort in daily life for the population of adults who

124    use CIs. A systematic review undertaken by the authors identified a number of PROMs

125    containing items considered to measure perceived listening effort in hearing loss. However,

126    no dedicated measures that have been developed with substantial input from CI recipients, or

127    validated in a large, representative sample, were found (Hughes et al. 2017a; Hughes et al.

128    2017b).

129         International guidance and consensus-based standards are available to guide PROM

130    development and validation (Mokkink et al. 2018; Patrick et al. 2011a; Patrick et al. 2011b).

131    The PROM development process includes content generation through direct engagement with

132    members of the target population; pretesting and qualitative refinement; and psychometric

133    evaluation to establish the new PROM's measurement properties (FDA 2009; FDA 2020).

134    Psychometric evaluation concerns the measurement of the target construct which is often

135    largely unobservable (e.g., effort, pain or fatigue). It aims to establish whether the

136    conceptualisation of the variable of interest (i.e., as defined by the PROM's content) has been

137    operationalised successfully (Hobart and Cano 2009). There are a number of psychometric

138    methods that may be used to evaluate a PROM's measurement characteristics, with Classical

139    Test Theory (CTT) and Item Response Theory (IRT) being two commonly used methods.

140    CTT is the most widely used of these methods and has dominated the field of psychometrics

141    for the last century (Streiner et al. 2015). Also known as weak true score theory, CTT

142    examines how measurement error affects rating scale scores. It proposes that an observed

143    score may be broken down into a true score and an error score, and although these values are

144    unknown, they may be estimated (Cano et al. 2011). CTT methods construct PROMs that

145    generate ordinal (rather than interval-level) data. The data are considered to represent ordered

146    counts, and the difference between adjacent categories on a response scale cannot be assumed

147    to have the same meaning (e.g., the distance between "sometimes" and "frequently" may not

148    be the same as "frequently" to "always") (Stevens 1946) . CTT evaluation is based on the

149    instrument as a whole and uses evidence primarily from correlations and descriptive statistics

150    (Hobart & Cano 2009).

151         Rasch measurement theory, a variant of IRT, is a modern psychometric method that is

152    used increasingly alongside CTT to develop and validate PROMs (Aryadoust et al. 2019).

153    Developed by Danish mathematician Georg Rasch, the Rasch model is a mathematical ideal

154    which specifies a set of criteria for the construction of interval level measures. The model

155    states that the probability that a person will affirm an item is a logistic function of the

156    difference between a person's trait level (this is expressed as person ability in Rasch terms)

157    and the amount of the trait expressed by the item (expressed as item difficulty in Rasch

158    terms), and is only a function of that difference (Rasch 1960). Therefore, the higher a

159    person's ability relative to item difficulty, the higher the probability of the person endorsing

160    that item. On this basis, a summed raw score may be used to derive an estimate of the target

161    trait (i.e., perceived listening effort) on an interval scale. Moreover, parameter separation, a

162    key characteristic of the Rasch model, means item difficulty is independent of the sample and

163    that a person's level of trait (i.e., ability) is independent of items (Bond & Fox 2015).

164         Rasch analysis is the term used to describe the formal evaluation of a PROM against

165    Rasch's mathematical measurement model. If data are found to conform to the model criteria,

166    then PROM developers can theoretically be confident that interval-level scales (measured

167    using the Rasch logit or log odds unit) have been constructed; that is, a respondent's answers

168    can be used to determine their precise location on a continuum that ranges from low to high

169    levels of perceived listening effort. Unlike CTT, the Rasch model enables estimates suitable

170    for individual-person analyses as well as group comparisons (Hobart & Cano 2009). These

171    two properties are of particular importance if a PROM's intended use is with individual

172    patients in the clinic. Because psychometric evaluation focuses on item-level rather than

173    scale-level assessment, a Rasch analysis also provides opportunity for PROM developers to

174    refine individual items and response scales, addressing potential sources of error and enabling

175    measurement precision to be improved (Boone 2016).

176        To satisfy the Rasch model, several criteria must be met. Firstly, an instrument's

177    items must demonstrate unidimensionality and local independence. For a scale to be

178    unidimensional, the items must only measure a single construct. Local independence means

179    items should only correlate through the latent trait that the test is measuring (i.e., a response

180    to an item must not influence responses to other items) (Lord et al. 1968). Thirdly, items must

181    fit the Rasch model expectations (i.e., item fit), as assessed through the examination of fit

182    statistics. Lastly, measures must be invariant such that the relative location of any two

183    persons on the construct continuum should be independent of the items used to make that

184    comparison (Hobart & Cano 2009).[1]

185        The opportunity for item refinement and construction of an instrument capable of

186    interval-level measurement means Rasch analysis could be argued to be a sensible choice for

187    undertaking a first evaluation of a new PROM's psychometric measurement properties.

188    Building upon developmental work reported elsewhere (Hughes et al. 2018), this paper

189    presents the quantitative refinement and initial validation of the Listening Effort

190    Questionnaire-Cochlear Implant (LEQ-CI), a new PROM measuring perceived listening

191    effort in daily life (Hughes et al. 2019). The specific study aims were: 1) to use Rasch

192    analysis to refine the LEQ-CI's rating scales and response categories; and 2) to assess the fit

193    of the LEQ-CI to the Rasch model. The overall goal is to make available to clinicians and

194    researchers a carefully developed and validated PROM of perceived listening effort that is

195    appropriate for use both in research and clinical practice.

196

---

[1] For readers interested in undertaking further study of psychometrics including Classical Test Theory and Rasch measurement theory, we suggest Bond et al. (2020), Boone et al. (2014), and Streiner et al. (2015) as introductory texts.

197

## MATERIALS AND METHODS

199    **Participants**

200         Linacre (1994) specifies a minimum sample size of 250 participants to increase

201    precision and robustness of estimates when conducting a Rasch analysis. Altogether, a total

202    of 511 cochlear implant patients from five regional National Health Service (NHS) auditory

203    implant centres in the UK were approached. Patients were invited to participate if they met

204    the following inclusion criteria: 1) they were adults aged 18 years or older; 2) they had a

205    diagnosis of severe-profound sensorineural hearing loss (SNHL); 3) they were either a CI

206    candidate or CI recipient according to the UK candidacy criteria (National Institute for Health

207    and Care Excellence (NICE) 2009); 4) they were able to read/write in competent English; 5)

208    they had no history of medical conditions that would preclude their ability to self-complete

209    the LEQ-CI; and 6) they were able to return the completed questionnaire by post. Both

210    candidates and recipients were invited to participate to increase the likelihood that persons

211    representing the full continuum of the latent trait (i.e., low to high levels of perceived

212    listening effort) were included in the study sample. Participating CI teams identified potential

213    participants and sent each participant an invitation letter, patient information sheet, and a

214    booklet comprised of a demographic questionnaire and the draft LEQ-CI[29]. A postage-paid

215    reply envelope was included for returning the completed booklet to the study team.

216

217    **Ethical considerations**

218         The study was approved by the North East - Newcastle and North Tyneside 2

219    Research Ethics Committee (Ref: 18/NE/0320) and the Swansea University-Swansea Bay

220    University Health Board's Joint Study Review Committee.

221

**The Listening Effort Questionnaire – Cochlear Implant (LEQ-CI)**

222

223     The LEQ-CI was developed to measure perceived listening effort in daily life in adult

224     CI patients. It is intended for use as both a candidacy and outcome measure in routine clinical

225     practice as well as in clinical research studies. The LEQ-CI's underlying conceptual

226     framework (see Figure 1) was developed from a mixed methods qualitative study that utilised

227     in-depth focus groups and a follow-up postal survey with cochlear implant patients (Hughes

228     et al. 2018). Findings from these first-hand accounts suggested that perceived listening effort

229     in daily life is a complex construct that includes the mental energy needed to attend to and

230     process an auditory signal while adapting and compensating for hearing loss. Effort is subject

231     to a cost-benefit analysis and is mediated by a number of motivational factors, particularly

232     that of social connectedness and pleasure. These findings also supported the heuristic

233     Framework for Understanding Effortful Listening (FUEL). The FUEL proposes that listening

234     effort depends not only on hearing difficulties and task demands, but also on the listener's

235     motivation to expend mental effort in the challenging situations of everyday life (Hughes et

236     al. 2017b; Pichora-Fuller et al. 2016).

237     An item pool was generated from the focus group accounts; a review of the literature;

238     and extant PROMs measuring relevant concepts within the conceptual framework (e.g.,

239     social connectedness, attention, and memory) (Hughes et al. 2018). Item candidates were

240     selected for inclusion in the draft LEQ-CI using the PROMIS® qualitative item review

241     procedures (DeWalt et al. 2007). The draft instrument (including respondents' instructions,

242     items and rating scales) was appraised for relevance, comprehensiveness, comprehensibility

243     and acceptability by an online survey of subject matter experts and also by CI patients in a

244     series of cognitive interviews (Hughes et al. 2019a). Qualitative findings from the online

245     survey and cognitive interviews were used to refine the LEQ-CI prior to undertaking further

246     refinement or performing an initial validation using Rasch analysis.

247          Qualitative refinement produced a draft LEQ-CI comprised of 29 items (LEQ-CI[29])

248    covering four domains (see Table 1): 1) the effort of attending (5 items); 2) the effort of

249    processing (9 items); 3) the effort associated with adapting and compensating for a hearing

250    loss (4 items); and 4) motivation (11 items). Items were measured using 5-point or 7-point

251    ordinal scales assessing either frequency (i.e., 1 = never; 2 = rarely; 3 = occasionally; 4 =

252    sometimes; 5 = frequently; 6 = usually; 7 = always) or intensity (i.e., 1 = not at all; 2 =

253    slightly; 3 = moderately; 4 = quite a bit; 5 = extremely), with higher scores indicating greater

254    levels of perceived listening effort. All relevant items were included in the draft instrument to

255    maximise coverage of the target construct (Bond & Fox 2015). A broad response scale was

256    utilised to enable greater differentiation in the judgements being made (Krosnick & Presser

257    2010) and to enable the use of Rasch analysis for optimisation of the rating scale response

258    categories (Bond & Fox 2015). The inclusion of a relatively large number of items reflected

259    the expectation that some items could be removed to optimise model fit without sacrificing

260    construct coverage. Readability of the full participant questionnaire booklet (including both

261    the LEQ-CI[29] and demographic questionnaire) was scored using the web-based application

262    Readable (www.readable.io). The booklet received a Flesch-Kincaid Grade Level score of

263    4.8 and SMOG Index of 7.8 (Flesch 1948; McLaughlin 1969). These results suggested the

264    wording of the instructions, items and response scales of the booklet would be understood by

265    a respondent with literacy skills equivalent to a 4th grade reading level based on the

266    American education system. A SMOG score 7.8 suggested the LEQ-CI[29] would be

267    understood by 93% of the UK population (NHS Digital Service Manual 2019). Estimated

268    reading time was approximately seven minutes.

269

270    **Data collection**

271     Participating CI clinical teams sent a copy of the study information leaflet, the LEQ-

272     CI[29], a demographic questionnaire, and a reply-paid envelope to CI patients who met the

273     study eligibility criteria. Patients who wished to participate completed the questionnaires and

274     returned these directly to the research team in the reply-paid envelope. Paper-based

275     questionnaires were used to optimise responding in an older CI patient population and to

276     promote inclusivity by limiting digital exclusion (Rowen et al. 2019; Smith et al. 2019). To

277     ensure participant anonymity, each booklet was coded with a unique identifier and no

278     identifiable information was collected. The clinical team at each CI centre held a master list

279     linking the identifiers to patients for purposes of data chasing (i.e., a reminder letter sent two

280     weeks after the initial mail out) and to pay each participant an honorarium in the form of a

281     £10 GBP retail gift card as a gesture of appreciation for their time and effort. All participant

282     follow-up was undertaken solely by the clinical team. To maintain anonymity, consent was

283     presumed when participants completed and returned the questionnaire booklet directly to the

284     research team (UK Data Service 2020).

285

286     **Missing data**

287     Any LEQ-CI[29] items not answered by respondents were considered missing data. The

288     full dataset, including cases with missing data, was included for analysis using Winsteps

289     analysis software (Version 4.4.7, www.winsteps.com). Unlike CTT which requires a

290     complete dataset, it is standard practice in Rasch analysis to include cases with missing data.

291     Winsteps's use of Joint Maximum Likelihood Estimation (JMLE) enables parameter

292     estimation for each missing case via a likelihood function based on the available data

293     (Linacre 2020a; Waterbury 2019).

294

295     **Data analysis plan**

296       Data analysis was undertaken in two stages: 1) a description of sample characteristics,

297    response rates, data quality and missing data; and 2) psychometric evaluation using Rasch

298    analysis, performed with Winsteps using JMLE and the partial credit model (Masters 1982).

299    The partial credit model is an unconstrained model suitable for polytomous data where the

300    response structure is permitted to vary across items (Bond and Fox 2015). The specific steps

301    of the Rasch analysis are described below.

302

303    <u>Assessment of rating scale functioning</u> - The rating scale structure of a PROM must fulfil

304    several criteria in order to satisfy the Rasch model requirements (Linacre 2002a). Firstly, all

305    items must orient with the latent variable. Any reverse-scored items should be re-scored, as

306    the rating scale categories may otherwise function differently. Secondly, a greater level of

307    perceived listening effort by a person should equate to higher scores on the LEQ-CI. As such,

308    it is expected that the higher the level of effort, the higher the response category that will be

309    endorsed, with response category thresholds advancing monotonically across the

310    measurement continuum. For each item, a category probability curve shows the probability of

311    observing each ordered category according to the Rasch model. If category thresholds do not

312    progress in a linear manner (i.e., along a continuum partitioned into equal, contiguous

313    intervals), or if the distance between two categories is judged to be inadequate, thresholds are

314    considered disordered. Disordering occurs when respondents select a response option that is

315    inconsistent with their ability and implies that a rating scale's categories may be confusing or

316    difficult to use (Boone et al. 2014; Pallant & Tennant 2007). Category disordering may be

317    resolved by collapsing one or more rating scale categories with adjacent categories (Boone &

318    Noltemeyer 2017). Thirdly, each response category should be endorsed by a minimum of ten

319    persons per response category as low category endorsement could mean the step calibration is

320    imprecisely estimated and potentially unstable (Linacre 2002a). A uniform distribution of

321    observations across response categories is optimal for step calibration. Lastly, category fit

322    statistics (i.e., as indices of how well the rating scale category structure meets the Rasch

323    model requirements) must fall within an acceptable range. Infit and outfit are the indices used

324    to assess model fit. Infit denotes inlier-sensitive or information-weighted fit, whilst outfit

325    refers to outlier-sensitive fit. Fit statistics are reported as mean square values (MNSQ), or as

326    standardised (ZSTD) values. MNSQ is the mean of the squared residuals for an item and is

327    sample independent, whilst ZSTD is a transformation of the mean square value with a sample

328    size correction. Outfit mean square values less than 2.0 logits (MNSQ < 2.0 logits) were

329    considered to be evidence of acceptable category fit (Bond et al. 2020).

330

331    Assessment of item and person fit to the model - Fit statistics show the extent that item

332    performance and person ability differ from the Rasch modelled expectations for fundamental

333    measurement. Mean-square values were reported using Winsteps, with values between 0.5

334    and 1.5 logits indicating acceptable fit (Linacre 2002b; Wright & Linacre 1994). Item fit was

335    confirmed by inspecting the item characteristic curves (ICCs) for each item. The ICC is an

336    ogive-shaped plot of the probabilities of responding to an item for any value of the

337    underlying trait (Bond & Fox 2015). The empirical and model ICCs should present a close

338    alignment of the observed and predicted scores for an item.

339          Person-fit measurement aims to identify individuals in the sample whose response

340    patterns deviate from the Rasch-modelled expectations (Boone et al. 2014). Person misfit has

341    the potential to affect item fit and, as a consequence, a scale's internal construct validity.

342    Person fit was assessed in the same manner as item fit. Examination of infit and outfit

343    statistics with mean square values outside the 0.5 – 1.5 logit range were considered evidence

344    of misfit (Linacre 2002b). Underfit, defined as high mean square values (MNSQ > 1.5), is

345    indicative of noise or randomness in the data with potential to influence item calibration;

346    therefore, it is standard practice to remove underfitting persons from the sample when

347    evaluating items' fit to the Rasch model (Bond & Fox 2015).

348

349    Assessment of unidimensionality and local independence - Unidimensionality (i.e., the set of

350    items under assessment represents a single construct) and local independence of items (i.e.,

351    the entire correlation between the items is captured by the latent trait) are core requirements

352    of the Rasch model (Hagquist et al. 2009). Unidimensionality was assessed in Winsteps using

353    a principal component analysis of the residuals (PCAR) (Linacre 2018). The following

354    criteria were used to determine whether the LEQ-CI met the requirement of

355    unidimensionality: 1) an eigenvalue less than 3.0 on the first residual contrast (Fan & Bond

356    2019; Linacre 2018); 2) the percentage variance explained by the first contrast of 5% or less

357    (Smith et al. 2007); and 3) disattenuated correlations between the person measures for the

358    item clusters on the first residual contrast greater than 0.70 (Linacre 2018). Local

359    independence was appraised by examining the item fit statistics for overfit and correlations

360    between items of the standardised residuals. Overfit and a substantial correlation of the

361    standardised residuals for two items was considered evidence of local dependence (Fan &

362    Bond 2019). A range of critical values for residual correlations have been reported in the

363    literature (Christensen et al. 2017). For this study, residual correlation values of $r < 0.4$ were

364    considered evidence of low dependency (Bond et al. 2020; Linacre 2020a).

365

366    Assessment of differential item functioning (DIF) - Differential item functioning (DIF)

367    analysis investigated whether the LEQ-CI's items measure perceived listening effort in daily

368    life in the same way for different sub-groups of the sample when both groups have equal

369    levels of effort. The presence of DIF implies a lack of measurement invariance that can

370    impact model fit (Tennant & Conaghan 2007). In the case of the LEQ-CI, DIF was explored

371    for gender (i.e., male and female) and age (i.e., adults aged less than 70 years and adults,

372    aged 70 years and older, as an approximation of the median age of the sample). Group

373    comparisons (e.g., males v females) were assessed for each item, examining probability of

374    DIF and effect size. Items found to have p-values less than 0.05 and an effect size (i.e., DIF

375    contrast value) of 0.64 or greater were considered to have sizable DIF (Linacre 2020).

376

377    Targeting of the scale - Targeting assesses the ability of an instrument's items to measure the

378    full range of persons in a sample. Items that are used to represent the latent construct (i.e.,

379    perceived listening effort) should collectively form a "difficulty" hierarchy, with difficulty

380    defined as the amount of trait represented by an item. The item hierarchy should range from

381    the item representing the least effort to the item measuring maximum effort. In a Rasch

382    analysis, item difficulty and person ability is measured using log-odds units or logits.[2] The

383    scale is always centred on zero logits, representing the item of average difficulty for the scale

384    (Tennant & Conaghan 2007). Item difficulties and person abilities are plotted on the same

385    logit scale and, for well-targeted instruments (i.e., one that is neither too easy nor too

386    difficult), the distribution of persons should closely match the distribution of items. The mean

387    person location and mean item location should correspond, with both location scores close to

388    zero. Targeting of the LEQ-CI was assessed by inspecting: 1) the person-item distribution

389    map 2) the summary statistics of item and person measures; and 3) the location range for

390    items and persons on the Rasch scale.

391

---

[2] The terms "item difficulty" and "person ability" are used in this manuscript to maintain consistency with Rasch terminology. Item difficulty refers to the level of perceived listening effort (i.e., the latent trait) that an item measures. Person ability is an estimate of the underlying trait or attribute being measured in an individual. In the case of the LEQ-CI, high ability is indicative of a low level of perceived listening effort (i.e., a person who reports listening to be relatively effortless).

392    Assessment of reliability - In Rasch analysis, reliability is defined as "reproducibility of

393    relative measure location", or the probability that persons or items with high measures

394    actually do have higher measures than items or persons with low measures (Wright &

395    Masters 1982). Reliability is sample size dependent, and measured in Winsteps using

396    separation indices. Separation indices indicate the number of distinct levels of functioning

397    that can be distinguished; ability in the case of persons; and difficulty in the case of items

398    (Duncan et al. 2003; Wright & Masters 1982). As a heuristic, a person separation index of

399    1.50 may be considered to represent an acceptable level of separation, whereas an index of

400    2.00 to represent a good level of separation and index of 3.00 represent to an excellent level

401    of separation (Wright & Masters 1982). A second statistic, the person reliability index,

402    indicates the replicability of person ordering if the same sample were given another set of

403    parallel items measuring the same construct. Higher values are indicative of better reliability

404    with values exceeding thresholds of 0.7, 0.8, and 0.9 indicating "acceptable", "good", and

405    "excellent" replicability (Bond & Fox 2015). Item separation and reliability is reported in the

406    same manner as person separation and reliability. The item separation index is used to verify

407    the item hierarchy (i.e., ordering of items according to difficulty) as confirmation of the

408    construct validity of an instrument (Wright & Masters 1982). Low item separation ($< 3$)

409    imply that the person sample is not large enough to confirm the item difficulty hierarchy

410    (Linacre 2011). High item reliability suggests item ordering would be replicated if the LEQ-

411    CI was given to a new but identical person sample (Wright & Masters 1982).

412

413                                                      **RESULTS**

414          From a total of 511 questionnaires distributed, 330 cochlear implant patients

415    completed the LEQ-CI[29], a response rate of 64.6%. The respondents represented a wide age

416    range (mean = 62.5 SD = 15.05, range = 21 – 89) that was negatively skewed with one third

417    of respondents aged 70 years or older (n = 125, 37.9%). There was a higher proportion of

418    females (n = 199, 60.3%) to males, and the majority of respondents used one CI only (n =

419    193, 58.5%) or one CI and a hearing aid (n = 106, 32.1%). No notable differences between

420    the candidate and recipient sub-samples were found for age, gender, or age at diagnosis of

421    hearing loss which suggested the subgroups were broadly similar in their demographic

422    characteristics. The sample characteristics were consistent with demographic trends for this

423    patient group (Amin et al. 2020). The demographic characteristics for the sample are reported

424    in Table 2.

425

426    **Data quality**

427        All data were extracted from the questionnaire booklet and entered manually into

428    REDCap (Version 8.10.2), clinical data management software hosted by Swansea Trials Unit,

429    Swansea University. Double entry was completed by the lead author (SEH) for 10% (n=33)

430    of the returned questionnaires. Eight data entry errors were detected (0.58%) across 1386

431    datapoints.

432        The complete dataset for all 330 participants was exported from REDCap to an Excel

433    spreadsheet and uploaded to Winsteps. Prior to upload into Winsteps, the Excel spreadsheet

434    was checked for formatting, and all reverse-scored items were transformed.

435

436    **Missing data**

437        Responses were missing for 73 items (0.007%) from the LEQ-CI[29] dataset of 9,570

438    (29 items x 330 respondents) responses. Eight respondents failed to complete one or more

439    pages of the LEQ-CI[29] accounting for 72.6% (n = 53) of the missing responses. On a per item

440    basis, the median number of missing responses was 2 (range = 0 – 6 missing responses,

441   corresponding to missing response rates of 0% - 1.8%). No further discernible response

442   patterns were identified from visual inspection of the data.

443

444   **Assessment of rating scale functioning and instrument refinement**

445            The first step of the Rasch analysis involved assessment of rating scale of functioning

446   for the LEQ-CI[29]'s items, following the guidelines proposed by Linacre (2002a). The

447   majority of the LEQ-CI[29]'s items showed multiple problems with category functioning when

448   5-point or 7-point rating scale structures were used. Categories failed to advance

449   monotonically for six items; disordered category fit statistics were reported for nine items;

450   and six items had response categories with fewer than 10 respondents. The category

451   probability curves showed threshold disordering and a lack modal peaks for all response

452   options. To improve rating scale function, adjacent response categories were collapsed to

453   construct either a 3-point or a 4-point rating scale. Restructuring of the rating scale categories

454   proceeded iteratively on a per item basis with the aim of optimising the response structure.

455   With each iteration, the revised rating scale was reviewed to check for uniform category

456   endorsement, to establish whether categories advanced monotonically and to confirm that

457   outfit MNSQ values were less than 2.0 logits as evidence of model fit. The category

458   probability curves were inspected, checking for modal peaks and advancing categories (Bond

459   & Fox 2015; Linacre 2002a; Wright 1996).

460            Despite refinements to the rating scale categories, eight items continued to exhibit

461   problems with scale functioning. No discernible patterns in terms of these items' content,

462   difficulty, or response category structure, were identified.  These items were removed to yield

463   a 21-item version with all items satisfying the essential criteria for rating scale functioning

464   (Linacre 2002a). The LEQ-CI[21]'s rating scales were found to be oriented with the latent

465   variable (i.e., increasing category score = increasing effort requirement), average measures

466    were found to advance monotonically, and outfit mean-square values were less than 2.0

467    logits. Table 3 presents a summary of the category structure for the refined rating scales of

468    the LEQ-CI[21]. Most category probability curves showed good rating scale functioning;

469    however, six items had category probability curves with ordered thresholds that failed to

470    show clear modal peaks for intermediate categories. The lack of modal peaks was attributed

471    to a lack of uniformity in category endorsement and was not considered to be detrimental to

472    rating scale functioning (see Figure, Supplemental Digital Content 1, showing the category

473    probability curves for the 29-item and 21-item versions of the LEQ-CI). Following rating

474    scale refinement, a further Rasch analysis was undertaken to establish whether the 21-item

475    version of the LEQ-CI met the Rasch model criteria for successful interval-level

476    measurement. The results of this analysis are reported below.

477

478    **Assessment of item and person fit to the model**

479        High, positive person fit residuals (MNSQ values > +2.0) are indicative of an

480    abnormal, random response pattern. Forty-six persons (13.9%) with MNSQ-values > +2.0

481    were identified as underfitting in the model and these persons were temporarily removed

482    from the dataset (Bond et al 2020) to ascertain their contribution to distortion of model fit.

483    Rasch analyses were conducted for the full data set (n = 330) and with underfitting persons

484    removed (n = 284). No notable differences between the analyses were identified; therefore,

485    the results of the Rasch analysis conducted with the full sample were reported (see Table,

486    Supplementary Digital Content 2, which presents the Rasch analyses for the full sample

487    (n=330) and with underfitting persons removed (n = 284)). All 21 items showed appropriate

488    fit to the model, with item fit statistics (mean-square values) falling within the 0.5 - 1.5 logits

489    range for productive measurement (see Table 4) (Linacre 2002b). The empirical range was

490    0.68 - 1.52. Empirical item characteristic curves (ICCs) showed good fit to the theoretical

491    ICCs with data points falling within the 95% confidence bands (see Figure, Supplemental

492    Digital Content 3, which shows the ICCs for the LEQ-CI[21]).

493

494    **Assessment of unidimensionality and local independence**

495            The LEQ-CI was found to meet the Rasch model requirements of unidimensionality

496    and local independence. The results of the principal component analysis of the residuals

497    (PCAR) supported the existence of only one Rasch dimension. The observed variance was

498    56.5%, and the unexplained variance of the first residual contrast was 5.1% (Eigenvalue =

499    2.47). Disattenuated correlations of item clusters (i.e., a group of items based on a cluster

500    analysis of the PCAR loadings) were above the recommended minimum of 0.7 (empirical

501    range 0.89 – 0.98), and suggested the clusters measure a single construct. A second PCAR

502    was performed using Winsteps-generated Rasch simulated data that was equivalent to the

503    empirical data file (i.e., in terms of person measures, item difficulties, and rating scale

504    structure). The simulated data were used to verify the empirical data, establishing whether the

505    observed eigenvalue of 2.47 was due to a second dimension or the result of an artefact (noise)

506    in the data structure. The unexplained variance for the Rasch-modelled data was 3.3%

507    (Eigenvalue = 1.34), which corresponded closely to the empirical data and provided further

508    evidence of unidimensionality. None of the items were overfitting and standardised residual

509    correlation values were less than 0.33 for the empirical dataset, indicative of local

510    independence.

511

512    **Assessment of differential item functioning (DIF)**

513            DIF analyses were undertaken for gender and age (see Tables, Supplemental Digital

514    Content 4, which presents the results of the DIF analyses for the LEQ-CI[21]). No evidence of

515    DIF was identified. Although two items were found to have p-values indicating statistical

516    significance (p < 0.05), the DIF contrast values (i.e. effect sizes) did not exceed the 0.64

517    threshold for DIF in terms of effect size (Linacre 2020b).

518

519    **Targeting of the scale**

520            The person-item distribution map showed good targeting of items to persons (see

521    Figure 2). The mean person ability was -0.32, suggesting persons had a slightly lower

522    requirement for effort relative to the mean item difficulty (which is always "0"). The full

523    sample ability ranged from -4.58 to 3.84 logits, with non-extreme person locations (n = 327)

524    ranging from -3.54 to 2.49 logits. To illustrate the meaning of the scale as a continuum of

525    perceived listening effort, Item 21 (e103) was located at a logit of 1.05 on the scale. At this

526    location, Item 21 provides content coverage for a person with moderately high level of

527    perceived listening effort (see Figure 2 and Table 5). Conversely. Item 3 (nq4) was located at

528    a logit of -0.81, providing content coverage for persons with a relatively low level of

529    perceived listening effort (see Figure 2 and Table 5). The person with the lowest level of

530    perceived listening effort was located at -4.58 logits, whilst the person with the highest level

531    of perceived listening effort was located at 3.84 logits. Overall, the items were found to

532    provide good coverage for the majority of the sample (item measure range -0.81 to 1.05);

533    however, a lack of items at the extreme ends of the effort continuum suggested that precision

534    may be reduced for the measurement of persons with extremely high or extremely low levels

535    of perceived listening effort.

536

537    **Assessment of reliability**

538            The LEQ-CI[21] person separation reliability was 0.91 and the person separation index

539    was 3.28, suggesting excellent reliability and an excellent level of separation distinguishing

540    approximately four levels of person ability defined as "low"; "below average"; "above

541    average"; and "high" levels of perceived listening effort (Duncan et al. 2003; Wright &

542    Masters 1982). The item separation index was 9.69, indicating the sample size was sufficient

543    to confirm the ordering of items on the listening effort continuum (i.e., construct validity).

544    **Raw score to Rasch score conversion**

545         The raw score for the LEQ-CI[21] was calculated by summing the item responses for all

546    21-items. Rasch scores were calculated in Winsteps and rescaled to 0-100 values as an index

547    of perceived listening effort (see Table, Supplemental Digital Content 5 which presents the

548    conversion table).

549                                   **DISCUSSION**

550         The present study is an essential step in developing a clinically useful and validated

551    measure of perceived listening effort in daily life. The LEQ-CI[21] is the first hearing-related

552    PROM to be developed specifically for the measurement of perceived listening effort in

553    adults with hearing loss, and one of only a few hearing-related PROMs to be developed and

554    validated using modern psychometric techniques. The results of this initial validation study

555    using Rasch analysis showed that the 21-item version of the LEQ-CI satisfies the Rasch

556    model requirements for interval-level measurement and, therefore, has the potential to be

557    used as a clinical tool for monitoring individual changes in perceived listening effort and as a

558    research tool for group comparisons using parametric statistical tests (Browne & Cano 2019).

559         In the context of listening effort assessment, using Rasch analysis to develop the

560    LEQ-CI offered several advantages. Firstly, the Rasch model creates a "fixed ruler" that

561    represents an effort continuum. The LEQ-CI's raw scores, when expressed as a summed total

562    in logits or converted to a 0-100 scale, may be used as an index of a person's perceived level

563    of listening effort (Wright & Masters 1982). This transformation renders the meaning of

564    LEQ-CI's total scores clearer and more easily interpretable. Secondly, performance outcome

565    measures of listening effort, such as tasks carried out under controlled conditions (e.g., dual

566    task paradigm), could be co-calibrated with the LEQ-CI[21] using the same interval scale, thus

567    enabling the performance-based measure to be interpreted using the content of the LEQ-CI

568    (Regnault et al. 2020). Co-calibration is of particular relevance to the study of listening effort

569    as the relationships among reported findings from behavioural, physiological and self-report

570    measures are, as yet, not well understood (Alhanbali et al. 2019; Strand et al. 2018).

571          Using Rasch analysis to refine and validate the LEQ-CI also facilitates its use across

572    different cultural contexts, meaning it could potentially be a suitable tool for international

573    research (Riff et al. 2017; Tennant et al. 2004). If translated into other languages or used

574    within different cultural contexts, evaluation of measurement invariance through DIF analysis

575    could enable adjustments to the LEQ-CI[21] (i.e., removal of items with DIF) to create a

576    measure that is invariant, thus ensuring the equivalence of any translated versions as well as

577    the cultural appropriateness of included items and their response scales. A culturally and

578    linguistic invariant measure supports the pooling of data, enabling CI research to be

579    undertaken with large datasets, a phenomenon that is, as yet, relatively uncommon in the CI

580    literature (Boisvert et al. 2020).

581          Lastly, the Rasch model would, in turn, support the development of short forms and

582    multiple versions of an instrument, due to psychometric analysis occurring at item rather than

583    at scale level. One limitation of CTT is its focus on an instrument's total score such that

584    evidence of validation relates to the instrument in its entirety, meaning it cannot be applied to

585    individual items (McKenna et al. 2019). Information on item fit; item difficulty; and item

586    discrimination (derived through application of the Rasch model) can enable a small set of

587    items to be selected that will optimise precision, whilst ensuring full coverage of persons

588    across the Rasch logit scale. The study findings, therefore, could support the use of the LEQ-

589    CI[21] in the development of individualised measures of perceived listening effort, such as

590    computer adaptive tests (CAT) where the aim is to minimise respondent burden without

591    sacrificing measurement accuracy (Kane et al. 2020; Petersen et al. 2018).

592          Evaluation of the LEQ-CI using Rasch analysis offers a theoretical contribution to the

593    conceptualisation of perceived listening effort in daily life by providing evidence in support

594    of the Framework for Understanding Effortful Listening (FUEL, Pichora-Fuller, 2016).

595    Building upon the existing body of work on mental effort and motivation, the FUEL proposes

596    motivational factors to have a modulating influence on effort deployment (Pichora-Fuller et

597    al. 2016). The finding that all of the domains represented in the LEQ-CI's conceptual

598    framework (including motivation) are components of a single dimension, defined as

599    perceived listening effort in daily life, contributes new evidence to the growing body of

600    literature on the involvement of motivational factors when listening is effortful (Koelewijn et

601    al. 2018; Picou & Ricketts 2014; Zekveld et al. 2019). More specifically, the LEQ-CI[21]'s

602    assessment of constructs such as pleasure (Matthen 2016); anxiety (Monzani et al. 2008);

603    social connectedness (Lee & Robbins 1995); and effort-reward imbalance (Siegrist 1996)

604    provides evidence in support of the contribution of social, emotional and psychological

605    constructs to perceived listening effort (Bruya & Tang 2018; Pichora-Fuller 2016).

606

607    **Limitations**

608          This study has begun to provide insights into the theoretical understanding of

609    perceived listening effort and its measurement. However, there are some limitations to the

610    study which must be considered when evaluating its potential contribution. Firstly,

611    refinement of the LEQ-CI[29], including item removal and the collapsing of rating scale

612    categories, was undertaken *a posteriori* based on empirical data, with fit to the Rasch model

613    confirmed for the refined 21-item instrument using the same dataset. To partly address this

614    limitation, simulated Rasch-modelled data based on the empirical dataset were used. Further

615  evaluation of the LEQ-CI[21] is now needed to verify fit to the Rasch model in a new sample of

616  adult cochlear implant patients and, in particular, CI candidates.

617          In addition, a change to the UK candidacy criteria for cochlear implantation was

618  implemented during data collection phase of the study, meaning patients who met the

619  updated criteria were not included in the study sample. Further research is needed to confirm

620  the fit of the LEQ-CI[21] to the Rasch model in a sample of adults meeting the revised UK CI

621  candidacy criteria, with particular attention given to whether the instrument is invariant in

622  this group when compared with the existing UK CI population.

623          Lastly, although targeting of items and person was good overall, incomplete item

624  coverage for persons at the extreme ends of the listening effort continuum represents another

625  limitation of the LEQ-CI[21]. Further work is now therefore needed to develop a better

626  understanding of how effort is perceived when an individual is listening at the limits of their

627  ability, under highly taxing acoustic conditions, or just prior to quitting a listening task

628  assessed as too arduous. Future work could also include a qualitative exploration of the

629  experiences of perceived listening effort from the perspective of adults with mild-moderate

630  hearing loss. These findings could then be used for purposes of  developing additional items

631  suitable for measuring persons with an extremely low or an extremely high level of perceived

632  listening effort. Further studies to explore the invariance of the LEQ-CI[21] across the full

633  range of hearing loss and different hearing devices is also needed.

634

635                                    **CONCLUSIONS**

636          The LEQ-CI[21] is the first PROM developed specifically to measure perceived

637  listening effort in daily life for adult cochlear implant candidates and recipients. It is one of

638  only a few hearing-specific PROMs to be developed using modern psychometric methods

639  and is the product of a body of work undertaken in line with international consensus-based

640    standards on PROM development. The current study began the task of establishing the LEQ-

641    CI[21] as a reliable and valid, interval-level measure of perceived listening effort in daily life

642    suitable for use in research and in clinical practice. The use of Rasch analysis helped to

643    engage the LEQ-CI[21] with theoretical frameworks in the published literature. Further

644    validation work is planned that will add to this body of evidence through assessment of the

645    LEQ-CI[21]'s construct validity, reliability, and responsiveness using traditional psychometric

646    approaches, as well as to study its implementation in clinical practice.(7218 words)

647

665

666     For further information about the LEQ-CI please contact: PROMS@swansea.ac.uk.

667

668                                              **REFERENCES**

669      Alhanbali, S., Dawes, P., Millman, R.E., et al. (2018). Measures of listening effort are
670              multidimensional. *Ear Hear*, 40(5), 1084-1097.

671      Amieva, H., Ouvrard, C., Meillon, C., et al. (2018). Death, depression, disability, and
672              dementia associated with self-reported hearing problems: A 25-year study. *J.
673              Gerontol*, 73, 1383–1389.

674      Amin, N., Wong, G., Nunn, T., et al. (2020). The Outcomes of Cochlear Implantation in
675              Elderly Patients: A Single United Kingdom Center Experience. *Ear Nose Throat J*,
676              0145561320910662.

677      Aryadoust, V., Tan, H.A.H., Ng, L.Y. (2019). A scientometric review of Rasch measurement:
678              The rise and progress of a specialty. *Front Psychol*, 10:2197.

679      Boisvert, I., Reis, M., Au, A., et al. (2020). Cochlear implantation outcomes in adults: A
680              scoping review. *PLoS One*, 15, e0232421.

681      Bond, T.G., Fox, C.M. (2015). *Applying the Rasch Model: Fundamental Measurement in the
682              Human Sciences* 3rd edition. London: Routledge.

683      Bond, T.G., Yan, Z., Heene, M. (2020). *Applying the Rasch Model: Fundamental
684              Measurement in the Human Sciences* 4th edition. London: Routledge.

685      Boone, W.J. (2016). Rasch analysis for instrument development: Why, when, and how? *CBE
686              Life Sci Educ*, 15:rm4.

687      Boone, W.J., Noltemeyer, A. (2017). Rasch analysis: A primer for school psychology
688              researchers and practitioners. *Cogent Educ*, 4, 1–13.

689      Boone, W.J., Yale, M.S., Staver, J.R. (2014). *Rasch Analysis in the Human Sciences.*
690              London: Springer.

691      Bräcker, T., Hellmiss, S., Batsoulis, C., et al. (2019). Introducing real-life listening features
692              into the clinical test environment: Part II: Measuring the hearing performance and
693              evaluating the listening effort of individuals with a hearing implant. *Cochlear
694              Implants Int*, 20, 165–175.

695      Browne, J.P., Cano, S.J. (2019). A Rasch measurement theory approach to improve the
696              interpretation of patient-reported outcomes. *Med Care*, 57, S18-S23.

697      Bruya, B., Tang, Y.Y. (2018). Is attention really effort? Revisiting Daniel Kahneman's
698              influential 1973 book *Attention and Effort. Front Psychol*, 9, 1133.

699   Cano, S.J., Hobart, J.C. (2011). The problem with health measurement. *Patient Prefer*
700           *Adherence*, 5, 279–290.

701   Christensen, K.B., Makransky, G., Horton, M. (2017). Critical values for Yen's Q3:
702           Identification of local dependence in the Rasch model using residual correlations.
703           *Appl Psychol Meas*, 41, 178–194.

704   Devlin, N.J., Appleby, J. (2010). *Getting the Most Out of Proms: Putting Health Outcomes at*
705           *the Heart of NHS Decision-making*, London: The King's Fund.

706   DeWalt, D.A., Rothrock, N., Yount, S., et al. (2007). Evaluation of item candidates: The
707           PROMIS qualitative item review. *Med Care*, 45, 12–21.

708   Duncan, P.W., Bode, R.K., Lai, S.M., et al. (2003). Rasch analysis of a new stroke-specific
709           outcome scale: The Stroke Impact Scale. *Arch Phys Med Rehabil*, 84, 950–963.

710   Edwards, B. (2016). A model of auditory-cognitive processing and relevance to clinical
711           applicability. *Ear Hear*, 37, 85S–91S.

712   Fan, J., Bond, T.G. (2019). Applying Rasch Measurement in Language Assessment:
713           Unidimensionality and Local Independence. In V. Aryadoust and M. Raquel, eds.
714           *Quantitative Data Analysis for Language Assessment Volume I: Fundamental*
715           *Techniques*. Abdingdon: Routledge.

716   FDA (2020). *Principles for Selecting, Developing, Modifying, and Adapting Patient-*
717           *Reported Outcome Instruments for Use in Medical Device Evaluation: Draft*
718           *Guidance for Industry and Food and Drug Administration Staff, And Other*
719           *Stakeholders*, Rockville, MD: U.S. Food and Drug Administration.

720   FDA (2009). Guidance for Industry Use in Medical Product Development to Support
721           Labelling Claims Guidance for Industry. *Clin Fed Regist*, 1–39.

722   Field, J., Holmes, M.M., Newell, D. (2019). PROMs data: can it be used to make decisions
723           for individual patients? A narrative review. *Patient Relat Outcome Meas*, 10, 233–
724           241.

725   Flesch, R. (1948). A new readability yardstick. *J Appl Psychol*, 32, 221–233.

726   Gagné, J., Besser, J., Lemke, U. (2017). Behavioral assessment of listening effort using a
727           dual-task paradigm: A review. *Trends Hear*, 21, 1–25.

728   Granberg, S., Dahlström, J., Möller, C., et al. (2014). The ICF Core Sets for hearing loss -
729           researcher perspective. Part I: Systematic review of outcome measures identified in
730           audiological research. *Int J Audiol*, 53, 65–76.

731   Hagquist, C., Bruce, M., Gustavsson, J.P. (2009). Using the Rasch model in nursing research:
732           An introduction and illustrative example. *Int J Nurs Stud*, 46, 380–393.

733   Hobart, J., Cano, S. (2009). Improving the evaluation of therapeutic interventions in multiple
734           sclerosis: The role of new psychometric methods. *Health Technol Assess*, 13, 1-177.

735    Holman, J.A., Drummond, A., Hughes, S.E., et al. (2019). Hearing impairment and daily-life
736            fatigue: A qualitative study. *Int J Audiol*, 58, 408–416.

737    Hornsby, B.W.Y. (2013). The effects of hearing aid use on listening effort and mental fatigue
738            associated with sustained speech processing demands. *Ear Hear*, 34, 523–534.

739    Hornsby, B.W.Y., Kipp, A.M. (2016). Subjective ratings of fatigue and vigor in adults with
740            hearing loss are driven by perceived hearing difficulties not degree of hearing loss.
741            *Ear Hear.*, 37, e1-10.

742    Hornsby, B.W.Y., Naylor, G., Bess, F.H. (2016). A taxonomy of fatigue concepts and their
743            relation to hearing loss. *Ear Hear*, 37, 136S-144S.

744    Hughes, S.E., Boisvert, I., Rapport, F., et al. (2019a). Measuring the experience of listening
745            effort in cochlear implant candidates and recipients: Development and pretesting of
746            the LEQ-CI. In *Poster presentation*. Southampton: BCIG Academic Meeting 2019.

747    Hughes, S.E., Hutchings, H.A., Dobbs, T.D., et al. (2017a). A systematic review and
748            narrative synthesis of the measurement properties of patient-reported outcome
749            measures (PROMs) used to assess listening effort in hearing loss. In *British Society of
750            Audiology*. Harrogate.

751    Hughes, S.E., Hutchings, H.A., Rapport, F., et al. (2017b). Qualitative data supporting the
752            FUEL: Perceived listening effort in cochlear implantation. In *Cognitive Hearing
753            Science for Communication*. Linköping, Sweden.

754    Hughes, S.E., Hutchings, H.A., Rapport, F.L., et al. (2018). Social connectedness and
755            perceived listening effort in adult cochlear implant users: A Grounded Theory to
756            establish content validity for a new patient reported outcome measure. *Ear Hear,* 39,
757            922-934.

758    Hughes, S.E., Rapport, F., Watkins, A., et al. (2019b). Study protocol for the validation of a
759            new patient-reported outcome measure (PROM) of listening effort in cochlear
760            implantation: The Listening Effort Questionnaire-Cochlear Implant (LEQ-CI). *BMJ
761            Open*, 9, 1–8.

762    Hughes, S.E., Rapport, F.L., Boisvert, I., et al. (2017c). Patient-reported outcome measures
763            (PROMs ) for assessing perceived listening effort in hearing loss: Protocol for a
764            systematic review. *BMJ Open*, 7:e014995.

765    Johnson, J., Xu, J., Cox, R., et al. (2015). A comparison of two methods for measuring
766            listening effort as part of an audiologic test battery. *Am J Audiol*, 24, 419–431.

767    Kane, L.T., Namdari, S., Plummer, O.R., et al. (2020). Use of computerized adaptive testing
768            to develop more concise patient-reported outcome measures. *JB JS Open Access*, 5,
769            e0052.

770    Koelewijn, T., Zekveld, A.A., Lunner, T., et al. (2018). The effect of reward on listening
771            effort as reflected by the pupil dilation response. *Hear Res*, 367, 106–112.

772   Krosnick, J. A., & Presser, S. (2010). Question and Questionnaire Design. In J.D. Wright and
773          P.V. Marsden, eds. *Handbook of Survey Research,* 2nd edition. San Diego, CA:
774          Elsevier.

775   Lee, R.M., Robbins, S.B. (1995). Measuring belongingness: The Social Connectedness and
776          the Social Assurance scales. *J Couns Psychol*, 42, 232–241.

777   Linacre, J.M. (1994). Sample size and item calibration or person measure stability. *Rasch
778          Meas Trans*, 7, 328.

779   Linacre, J. (2002a). Understanding Rasch measurement: Optimizing rating scale category
780          effectiveness. *J Appl Meas*, 3, 85–106.

781   Linacre, J.M. (2002b). What do infit and outfit, mean-square and standardized mean. *Rasch
782          Meas Trans*, 16, 878.

783   Linacre, J.M. (2018). Detecting multidimensionality in Rasch data using Winsteps Table 23.
784          Available at: https://youtu.be/sna19QemE50 [Accessed November 29, 2020].

785   Linacre, J.M. (2020a). *A User's Guide to WINSTEPS MINISTEPS Rasch-Model Computer
786          Programs*, Available at:
787          http://homes.jcu.edu.au/$\sim$edtgb/%5Cnpapers3://publication/uuid/D56B724A-
788          62FF-4D00-84E1-ECC888298B70.

789   Linacre, J.M. (2020b). Table 30.1 Differential item functioning DIR pairwise. Available at:
790          https://www.winsteps.com/winman/table30_1.htm [Accessed November 29, 2020]..

791   Lord, F.M., Novick, M.R., Birnbaum, A. (1968). *Statistical Theories of Mental Test Scores.*
792          Oxford: Addison-Wesley.

793   Masters, G.N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47, 149–174.

794   Matthen, M. (2016). Effort and displeasure in people who are hard of hearing. *Ear Hear*, 37,
795          28S–34S.

796   McGarrigle, R., Munro, K.J., Dawes, P., et al. (2014). Listening effort and fatigue: What
797          exactly are we measuring? A British Society of Audiology Cognition in Hearing
798          Special Interest Group "white paper". *Int J Audiol*, 53, 433–40.

799   McLaughlin, G.H. (1969). SMOG grading: A new readability formula. *J Read*, 12, 639–646.

800   McKenna, S.P., Heaney, A., Wilburn, J. et al. (2019). Measurement of patient-reported
801          outcomes. 1: The search for the Holy Grail. *J Med Econ*, 22, 516-522.
802
803   McRackan, T.R., Velozo, C.A., Holcomb, M.A., et al. (2017). Use of adult patient focus
804          groups to develop the initial item bank for a cochlear implant quality-of-life
805          instrument. *JAMA Otolaryngol Head Neck Surg*, 143, 975–982.

806   Miles, K., McMahon, C., Boisvert, I., et al. (2017). Objective assessment of listening effort:
807          coregistration of pupillometry and EEG. *Trends Hear*, 21, 1-13.

808   Mokkink, L.B., Prinsen, C.A.C., Patrick, D.L., et al. (2018). *COSMIN methodology for*
809          *systematic reviews of Patient - Reported Outcome Measures ( PROMs )*,

810   Monzani, D., Galeazzi, G., Genovese, E., et al. (2008). Psychological profile and social
811          behaviour of working adults with mild or moderate hearing loss. *Acta*
812          *Otorhinolaryngol Ital*, 28, 61–66.

813   Nachtegaal, J., Kuik, D.J., Anema, J.R., et al. (2009). Hearing status, need for recovery after
814          work, and psychosocial work characteristics: Results from an internet-based national
815          survey on hearing. *Int J Audiol*, 48, 684–91.

816   National Institute for Health and Clinical Excellence (NICE) (2009). *NICE technology*
817          *appraisal guidance 166: Cochlear implants for children and adults with severe to*
818          *profound deafness*. Available at: guidance.nice.org.uk/TA166/pdf/English [Accessed
819          March 24, 2020].

820   NHS Digital Service Manual (2019). Use a readability tool to prioritise content - NHS digital
821          service manual. Available at: https://service-manual.nhs.uk [Accessed July 22, 2020].

822   Ohlenforst, B., Zekveld, A. (2017). Effects of hearing impairment and hearing aid
823          amplification on listening effort: A systematic review. *Ear Hear,* 38, 267-281.

824   Pallant, J.F., & Tennant, A. (2007). An introduction to the Rasch measurement model: An
825          example using the Hospital Anxiety and Depression Scale (HADS). *Br J Clin*
826          *Psychol*, 46, 1–18.

827   Patrick, D.L., Burke, L.B., Gwaltney, C.J., et al. (2011a). Content validity - Establishing and
828          reporting the evidence in newly developed patient-reported outcomes (PRO)
829          instruments for medical product evaluation: ISPOR PRO good research practices task
830          force report: Part 2 - Assessing respondent understanding. *Value Health*, 14, 978–988.

831   Patrick, D.L., Burke, L.B., Gwaltney, C.J., et al. (2011b). Content validity–establishing and
832          reporting the evidence in newly developed patient-reported outcomes (PRO)
833          instruments for medical product evaluation: ISPOR PRO good research practices task
834          force report: part 1–eliciting concepts for a new PRO instrument. *Value Health*, 14,
835          967–77.

836   Petersen, M.A., Aaronson, N.K., Arraras, J.I., et al. (2018). The EORTC CAT Core—The
837          computer adaptive version of the EORTC QLQ-C30 questionnaire. *Eur J Cancer*,
838          100, 8–16.

839   Pichora-Fuller, M.K. (2016). How social psychological factors may modulate auditory and
840          cognitive functioning during listening. *Ear Hear*, 37, 92S–100S.

841   Pichora-Fuller, M.K., Kramer, S.E., Eckert, M.A., et al. (2016). Hearing impairment and
842          cognitive energy. *Ear Hear*, 37, 5S–27S.

843   Picou, E.M., Ricketts, T. a (2014). Increasing motivation changes subjective reports of
844          listening effort and choice of coping strategy. *Int J Audiol*, 53, 1–9.

845    Ramage-Morin, P.L. (2016). *Hearing Difficulties and Feelings of Social Isolation Among*
846            *Canadians Aged 45 or Older*, Available at: http://www.statcan.gc.ca/pub/82-003-
847            x/2016011/article/14671-eng.pdf [Accessed April 12, 2020].

848    Rasch, G. (1960). *Probabilistic Models for Some Intelligence and Attainment Tests*,
849            Copenhagen: Nielson and Lydiche.

850    Regnault, A., Meunier, J., Ciesluk, A., et al. (2020). Providing meaningful interpretation of a
851            performance outcome measure by co-calibration with a patient-reported outcome:
852            illustration with widely used multiple sclerosis measures. In Virtual ISOQOL 2020.

853    Riff, K.W.Y.W., Tsangaris, E., Goodacre, T., et al. (2017). International multiphase mixed
854            methods study protocol to develop a cross-cultural patient-reported outcome
855            instrument for children and young adults with cleft lip and/or palate (CLEFT-Q). *BMJ*
856            *Open*, 7, e015467.

857    Rowen, D., Carlton, J., Elliott, J. (2019). PROM validation using paper-based or online
858            surveys: Data collection methods affect the sociodemographic and health profile of
859            the sample. *Value Health*, 22, 845–850.

860    Siegrist, J. (1996). Adverse health effects of high-effort / low-reward conditions. *J Occup*
861            *Health Psychol*, 1, 27–41.

862    Smith, A.B., Wright, P., Selby, P.J., et al. (2007). A Rasch and factor analysis of the
863            Functional Assessment of Cancer Therapy-General (FACT-G). *Health Qual. Life*
864            *Outcomes*, 5, 1–10.

865    Smith, M.G., Witte, M., Rocha, S., et al. (2019). Effectiveness of incentives and follow-up on
866            increasing survey response rates and participation in field studies. *BMC Med Res*
867            *Methodol*, 19, 230.

868    Stevens, S.S. (1946). On the Theory of Scales of Measurement. *Science*, 103, 677–680.

869    Strand, J.F., Brown, V.A., Merchant, M.B., et al. (2018). Measuring listening effort:
870            convergent validity, sensitivity, and links with cognitive and personality measures. *J*
871            *Speech Lang Hear Res*, 61, 1463–1486.

872    Strand, J.F., Dillman-Hasso, N., Villanueva, J., et al. (2020). Understanding speech amid the
873            jingle and jangle: Recommendations for improving measurement practices in listening
874            effort research. *Pre-Print*. Available at: https://doi.org/10.31234/osf.io/3e7mf
875            [Accessed November 29, 2020].
876

877    Streiner, D.L., Norman, G.R., Cairney, J. (2015). *Health Measurement Scales: A Practical*
878            *Guide to Their Development and Use*, 5th Edition. Oxford: Oxford University Press.

879    Tennant, A., & Conaghan, P.G. (2007). The Rasch measurement model in rheumatology:
880            What is it and why use it? When should it be applied, and what should one look for in
881            a Rasch paper? *Arthritis Rheum*, 57, 1358-1362.

882    Tennant, A., Penta, M., Tesio, L., et al. (2004). Assessing and Adjusting for Cross-Cultural
883            Validity of Impairment and Activity Limitation Scales Through Differential Item

884           Functioning Within the Framework of the Rasch Model: The PRO-ESOR Project.
885           *Med Care*, 42, I37-I38.

886   UK Data Service (2020). Consent for Data Sharing. Available at:
887           https://ukdataservice.ac.uk/manage-data/legal-ethical/consent-data-
888           sharing/surveys.aspx [Accessed November 5, 2020].

889   Waterbury, G.T. (2019). Missing Data and the Rasch Model: The Effects of Missing Data
890           Mechanisms on Item Parameter Estimation. *J Appl Meas*, 20, 154–166.

891   Wright, B.D. (1996). Comparing Rasch measurement and factor analysis. *Struct Equ Model*,
892           3, 3–24.

893   Wright, B.D., Masters, N.G. (1982). *Rating Scale Analysis*, Chicago: MESA Press.

894   Wright, B.D., Stone, M.H. (1979). *Best Test Design*, Chicago: MESA Press.

895   Wright, Linacre (1994). Reasonable mean-square fit values. *Rasch Meas Trans*, 8, 370.

896   Zekveld, A.A., Koelewijn, T., Kramer, S.E. (2018). The pupil dilation response to auditory
897           stimuli: Current state of knowledge. *Trends Hear*, 22, 1-25.

898   Zekveld, A.A., van Scheepen, J.A.M., Versfeld, N.J., et al. (2019). Please try harder! The
899           influence of hearing status and evaluative feedback during listening on the pupil
900           dilation response, saliva-cortisol and saliva alpha-amylase levels. *Hear Res*, 381,
901           107768.

902

**FIGURE LEGEND**

904 Figure 1. The conceptual framework for measurement of perceived listening effort in daily

905 life for the adult cochlear implant population. This diagrammatic representation shows the

906 domains and concepts, as well as the expected relationships between concepts (i.e., items),

907 measured by the LEQ-CI[21]

908



909

910     Figure 2. The person-item map for the LEQ-CI[21] (n = 330) showing the distribution of

911     persons (upper half) and items (lower half). The x-axis displays the perceived listening effort

912     continuum on a logit scale and the y-axis displays frequency counts.

913



914

915    Table 1: Item stems per domain for the LEQ-CI showing all 29 items included in the draft instrument. An asterisk (*) denotes an item
916    removed due to misfitting the Rasch model. Remaining items comprise the LEQ-CI[21].
917

| Item No. | Item Label | Domain | Item stem |
|---|---|---|---|
| 1. | e398 | Attending | Have to stay alert |
| 2. | e319 | Attending | Strain to hear sounds around you |
| 3. | nq4 | Attending | Strain to hear speech |
| 4. | e320 | Attending | Figure out which sound to focus on |
| 5. | e31 | Attending | Listen as long as needed* |
| 6. | n2- | Processing | Anticipate what someone says* |
| 7. | e35 | Adapting & compensating | Ask others to repeat |
| 8. | n14 | Processing | Figure out what said - odd word or phrase* |
| 9. | e50 | Processing | Forget what other person just said |
| 10. | n16 | Processing | Take longer to understand |
| 11. | e323 | Processing | Able to listen to someone talk while doing something else |
| 12. | e71 | Processing | Listen and plan reply* |
| 13. | e365 | Attending | Listening just seems to happen* |
| 14. | e310 | Motivation | Stop listening because too much effort |
| 15. | e34 | Processing | Understand group talking - no background noise* |
| 16. | nq3 | Processing | Understand group talking - background noise |
| 17. | e121 | Adapting & compensating | Do things to make listener easier* |
| 18. | n11 | Adapting & compensating | Make others feel at ease |
| 19. | e344 | Motivation | Effort bother you |
| 20. | e213 | Motivation | Avoid situations because of effort needed to listen |
| 21. | e103 | Motivation | Give up on listening |
| 22. | e335 | Adapting & compensating | Plan day around effort* |
| 23. | e105 | Adapting & compensating | Run out of energy for listening |
| 24. | e214 | Adapting & compensating | Can be yourself |
| 25. | e118 | Motivation | Feel tense |
| 26. | e112 | Motivation | Choose to be alone because of listening effort |

| Item No. | Item Label | Domain | Item stem |
|---|---|---|---|
| 27. | e326 | Motivation | Stop doing things want to do |
| 28. | e363 | Motivation | Find pleasure in listening |
| 29. | e329 | Motivation | Feel connected with others |

918   *Item removed as misfitting the Rasch model

Table 2: Participants' demographic characteristics

|  | Total (N = 330) |
|---|---|
| **Patient age (years)** |  |
| Mean (SD) | 62.50 (15.1) |
| Median (range) | 66 (21 – 89) |
| Missing | 0 |
| **Gender** |  |
| Male | 131 (39.7%) |
| Female | 199 (60.3%) |
| Prefer to self-describe | 0 (0%) |
| Prefer not to say | 0 (0%) |
| Missing | 0 (0%) |
| **Employment status** |  |
| Employed, working full time | 76 (23.0%) |
| Employed, working part-time | 40 (12.1%) |
| Not employed | 11 (33.3%) |
| Retired | 177 (53.6%) |
| Full-time homemaker | 13 (3.9%) |
| Volunteer worker | 2 (0.6%) |
| Student | 2 (0.6%) |
| Other | 8 (2.4%) |
| Unknown | 1 (0.3%) |
| **Highest level of education** |  |
| Not applicable | 8 (2.4%) |
| Primary school | 1 (0.3%) |
| Secondary school | 96 (29.1%) |
| GCSE/O-level qualification | 52 (15.8%) |
| A-level/BTEC or equivalent post-16 education | 70 (21.2%) |
| Undergraduate degree | 21 (6.4%) |
| Graduate degree | 50 (15.2%) |
| Doctorate degree | 9 (2.7%) |
| Other | 20 (6.1%) |
| Unknown | 3 (0.9%) |
| **UK Region** |  |
| London | 71 (21.5%) |
| Midlands | 62 (18.8%) |
| Northwest | 97 (29.4%) |
| Scotland | 91 (27.6%) |
| Wales[*] | 9 (2.72%) |
| **Age at HL onset** |  |
| Mean (SD) | 28.6 (20.1) |
| Median (range) | 28.0 (0 – 78.0) |

| | |
|---|---|
| Missing | 10 (3.0%) |
| Age at receipt of first hearing device | |
|     Mean (SD) | 35.0 (20.9) |
|     Median (range) | 39.5 (0 – 84.0) |
|     Missing | 12 (0.04%) |
| Hearing devices used | |
|     One hearing aid | 7 (2.1%) |
|     Two hearing aids | 14 (4.2%) |
|     One cochlear implant | 193 (58.5%) |
|     One cochlear implant + one hearing aid | 106 (32.1%) |
|     Two cochlear implants | 7 (2.1%) |
|     None | 3 (0.9 %) |
|     Other | 0 (0%) |
|     Missing | 0 (0%) |
| Hours of hearing device use | |
|     Mean (SD) | 14.1 (2.9) |
|     Median (range) | 15.0 (0 – 24.0) |
|     Missing | 5 (1.5%) |
| Years implanted | |
|     Mean (SD) | 3.8 (4.8) |
|     Median (Range) | 2.0 (0 – 26.0) |
|     Missing | 26 (1.0%) |

*Fewer patients from the three Welsh CI centres were eligible to participate due to previous involvement as participants in LEQ-CI content validity studies. Nine responses were received from the 11 Welsh patients eligible to participate, a response rate of 81.8%.

Table 3: Summary of the category fit statistics for each of the LEQ-CI[21] rating scales. Category values are the raw category values that correspond to the category label, average measures can be seen to advance montonically, Outfit MnSq values < 2.0 are indicative of model fit.

| Item Group | Category Value | Average Measure (logits) | OUTFIT MnSq (logits) | Andrich Thresholds (logits) | Category Label |
|---|---|---|---|---|---|
| 1[a] | 1 | -1.25 | 1.13 | None | Never |
| | 3 | -0.47 | 0.87 | -1.96 | Occasionally |
| | 5 | 0.14 | 0.94 | -0.05 | Frequently |
| | 7 | 0.57 | 1.54 | 2.01 | Always |
| 2[b] | 1 | -0.67 | 1.12 | None | Not at all |
| | 2 | -0.27 | 0.68 | -0.77 | Slightly |
| | 3 | 0.22 | 0.75 | -0.09 | Moderately |
| | 5 | 0.74 | 1.06 | 0.86 | Extremely |
| 3[c] | 1 | -0.74 | 0.89 | None | Never |
| | 3 | 0.01 | 0.44 | -0.18 | 1 - 4 days |
| | 5 | 0.42 | 0.97 | 0.18 | 5 - 7 days |
| 4[d] | 1 | -0.99 | 0.90 | None | None of the time |
| | 3 | -0.19 | 0.96 | -1.80 | Some of the time |
| | 4 | 0.33 | 0.95 | 1.80 | Most of the time |
| 5[e] | 1 | -1.16 | 0.98 | None | Not at all tense |
| | 3 | -0.44 | 0.77 | -1.82 | Moderately tense |
| | 5 | 0.36 | 0.80 | 1.82 | Extremely tense |

Items are labelled according to their entry number in Winsteps and are grouped according to their rating scale structure:
[a]Items = Item 1, 4, 7, 9-11, 14, 16, 20, 21, 26, 27
[b]Items = Item 2, 3, 18, 19, 28, 29
[c]Items = 23
[d]Items = Item 24
[e]Items = Item 25

Table 4: Item fit statistics showing final item locations for the LEQ-CI[21] (n = 330)

| Item ID | Measure | Infit Statistics | | Outfit Statistics | |
|---|---|---|---|---|---|
| | | MNSQ | ZSTD | MNSQ | ZSTD |
| e103 | 1.05 | 1.03 | 0.43 | 0.90 | -0.55 |
| e326 | 0.84 | 0.97 | -0.31 | 0.86 | -0.97 |
| e112 | 0.64 | 0.82 | -2.32 | 0.72 | -2.32 |
| e50 | 0.63 | 1.14 | 1.65 | 1.08 | 0.65 |
| e310 | 0.60 | 0.87 | -1.72 | 0.86 | -1.18 |
| e213 | 0.48 | 1.03 | 0.39 | 0.93 | -0.66 |
| e105 | 0.32 | 0.91 | -1.11 | 0.7 | -1.08 |
| e320 | 0.09 | 1.18 | 2.27 | 1.21 | 2.18 |
| e329 | -0.02 | 0.73 | -3.87 | 0.81 | -2.19 |
| n16 | -0.11 | 0.80 | -2.79 | 0.93 | -0.81 |
| e363 | -0.14 | 1.02 | 0.27 | 1.10 | 1.04 |
| e214 | -0.20 | 1.03 | 0.41 | 1.02 | 0.19 |
| e319 | -0.23 | 1.22 | 2.84 | 1.25 | 2.41 |
| e35 | -0.25 | 0.66 | -5.12 | 0.7 | -3.33 |
| e323 | -0.30 | 1.43 | 5.07 | 1.47 | 5.01 |
| e398 | -0.43 | 1.48 | 5.61 | 1.52 | 5.51 |
| n11 | -0.45 | 1.27 | 3.41 | 1.37 | 3.28 |
| e344 | -0.47 | 0.72 | -4.21 | 0.68 | -3.50 |
| nq3 | -0.61 | 1.38 | 4.55 | 1.36 | 3.94 |
| e118 | -0.67 | 0.91 | -1.25 | 0.88 | -1.40 |
| nq4 | -0.81 | 0.70 | -4.34 | 0.83 | -1.30 |
| Mean | 0.00 | 1.01 | 0.0 | 1.01 | 0.20 |
| SD | 0.52 | 0.24 | 3.1 | 0.25 | 2.5 |

The MNSQ acceptable range for productive measurement = 0.5 – 1.5 logits.
MNSQ = mean square
ZSTD = Z-standardised statistics

1    Table 5: Extreme scores for items and persons from the sample (n = 330).

| Extreme Items/Persons | Total Score | Measure | Meaning of Item/Person Location |
|---|---|---|---|
| Item 21 (e103) | n/a | 1.05 | Low trait level (more likely to be endorsed by persons with moderately high levels of perceived listening effort) |
| Item 3 (nq4) | n/a | -0.81 | High trait level (more likely to be endorsed by persons with low levels of perceived listening effort) |
| Person 276 10/05/2021 11:26:00 | 124 | 3.84 | Person with highest level of perceived listening effort |
| Person 35 | 72 | 0.01 | Person with average levels of perceived listening effort |
| Person 203 | 21 | -4.58 | Person with the lowest level of perceived listening effort |

2
3

# LIST OF ABBREVIATIONS

| | |
|---|---|
| CAT | Computer adaptive test |
| CI | Cochlear implant |
| CTT | Classical Test Theory |
| DIF | Differential item functioning |
| EEG | Electroencephalography |
| FUEL | Framework for Understanding Effortful Listening |
| HA | Hearing aid |
| ICC | Item characteristic curve |
| IRT | Item response theory |
| JMLE | Joint Maximum Likelihood Estimation |
| LEQ-CI | Listening Effort Questionnaire – Cochlear Implant |
| MMHL | Mild-moderate hearing loss |
| MNSQ | Mean square |
| NHS | National Health Service |
| NICE | National Institute for Health and Care Excellence |
| PCAR | Principal component analysis of the residuals |
| PROM | Patient-reported outcome measure |
| PROMIS® | Patient-Reported Outcomes Measurement Information System |
| PSI | Person separation index |
| RMT | Rasch measurement theory |
| SNHL | Sensorineural hearing loss |
| UK | United Kingdom |
| VAS | Visual analogue scales |
| ZSTD | z-standardised |

7        **LIST OF SUPPLEMENTAL DIGITAL CONTENT**

8     Supplemental Digital Content 1. Category Probability Curves for the draft 29-item version

9     and the refined 21-item version of the LEQ-CI. docx

10

11     Supplemental Digital Content 2. Table comparing the findings of Rasch analyses undertaken

12     with the full sample (n = 330) and with underfitting persons removed (n = 284). docx

13

14     Supplemental Digital Content 3. Item Characteristic Curves for all items of the LEQ-CI[21].

15     pptx

16

17     Supplemental Digital Content 4. Tables presenting DIF analyses for gender (males v females)

18     and age (adults < 70 years and adults aged 70 years or older). docx

19

20     Supplemental Digital Content 5. Raw score to 0-100 Rasch score conversion table for use

21     with the LEQ-CI[21]. docx

# Supplemental Digital Content 1.

Category Probability Curves for the LEQ -CI[29] and LEQ-CI[21]

Rasch model category probability curves for each item of the original 29-item and the revised 21-item versions of the LEQ-CI. The x-axis represents the attribute in logits. The y-axis represents the probability of a response category being selected. The curves represent the likelihood that a respondent with a particular amount of the latent trait will select a category. The y-axis represents the expected probability of endorsement of any given category when a person responds to the item. The x-axis represents the person ability relative to the item difficulty, with origins set to 0. The scale is measured in logits (log odds units). Distances to the right of zero indicate higher and higher levels of effort. Locations to the left of zero indicate lower and lower levels of effort. Categories should advance in ascending order from left to right along the x-axis. Categories should each have a distinct peak indicating that it is the most probable (modal) category at that point on the latent variable. The crossover points between two curves are the equal-probability points or thresholds. Thresholds advance in line with categories, thresholds that fail to advance monotonically are considered disordered. The initial item category structure shows threshold disordering and a lack of modal peaks. The category probability curves have not been displayed for any items that have been removed from the LEQ-CI due to model misfit.

22

### LEQ-CI[29]
### LEQ-CI[21]



23

LEQ-CI[29]                                              LEQ-CI[21]



24

LEQ-CI[29]                                              LEQ-CI[21]



25

LEQ-CI[29]                                          LEQ-CI[21]



26

LEQ-CI[29]                                          LEQ-CI[21]



- Item removed

27

LEQ-CI[29]                                        LEQ-CI[21]

6. n_2



• Item removed

28

LEQ-CI[29]                                        LEQ-CI[21]

7. e35                                            7. e35



29

LEQ-CI[29]                                                           LEQ-CI[21]



- Item removed

LEQ-CI[29]                                                           LEQ-CI[21]

                                 

30

31

LEQ-CI[29]                                                         LEQ-CI[21]



32

LEQ-CI[29]                                                         LEQ-CI[21]



33

LEQ-CI[29]                                    LEQ-CI[21]



34

LEQ-CI[29]                                    LEQ-CI[21]



• Item removed

35

LEQ-CI[29]                                              LEQ-CI[21]



14. e310

• Item removed

36

LEQ-CI[29]                                              LEQ-CI[21]



15. e34r

• Item removed

37

LEQ-CI[29]                                                    LEQ-CI[21]



38

LEQ-CI[29]                                                    LEQ-CI[21]

- Item removed



39

LEQ-CI[29]                                            LEQ-CI[21]



18. n11                                               18. n11

40

LEQ-CI[29]                                            LEQ-CI[21]



19. e344                                              19. e344

41

LEQ-CI[29]                                              LEQ-CI[21]



42

LEQ-CI[29]                                              LEQ-CI[21]



43

LEQ-CI[29]                                    LEQ-CI[21]



44

LEQ-CI[29]                                    LEQ-CI[21]



45

LEQ-CI[29]                                              LEQ-CI[21]



46

LEQ-CI[29]                                              LEQ-CI[21]



47

LEQ-CI[29]                                    LEQ-CI[21]



48

LEQ-CI[29]                                    LEQ-CI[21]



49

LEQ-CI[29]                                              LEQ-CI[21]

28. e363r                                               28. e363r

50

LEQ-CI[29]                                              LEQ-CI[21]

29. e329r                                               29. e329r

51

Supplemental Digital Content 2: Table comparing the Rasch analyses undertaken with the full sample (n = 330) and with underfitting persons removed (n = 284)

|  | Full sample (n = 330) | Underfitting persons removed (n = 284) |
|---|---|---|
| Rasch variance explained by measures | 56.5% | 60.8% |
| Eigenvalue for PCA of 1$^{st}$ residual | 2.48 | 2.33 |
| Range disattenuated correlations of 1$^{st}$ residual | 0.89 - 0.98 | 0.92 – 1.00 |
| Mean person location | -0.31 | -0.39 |
| Person separation | 3.28 | 3.62 |
| Person reliability | 0.91 | 0.94 |
| Item separation | 9.69 | 10.07 |

# Supplemental Digital Content 3.

### Empirical + theoretical item characteristic curves (ICCs) for the LEQ -CI[21]
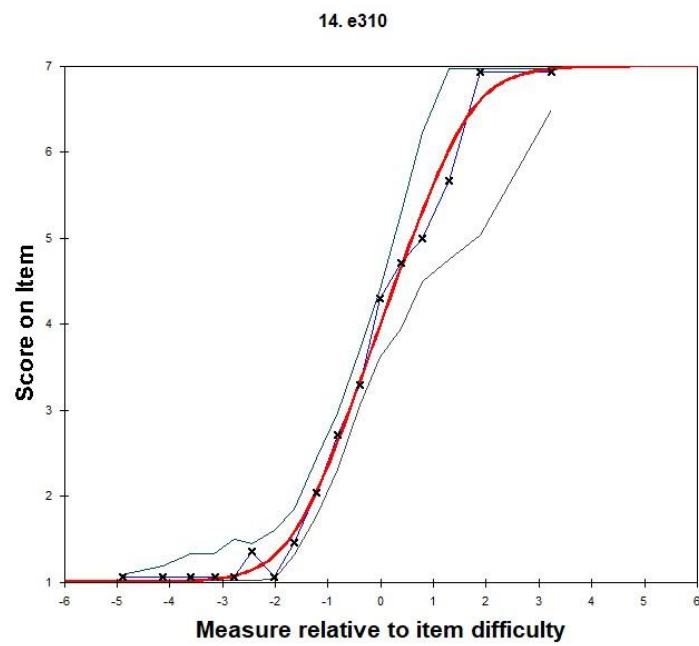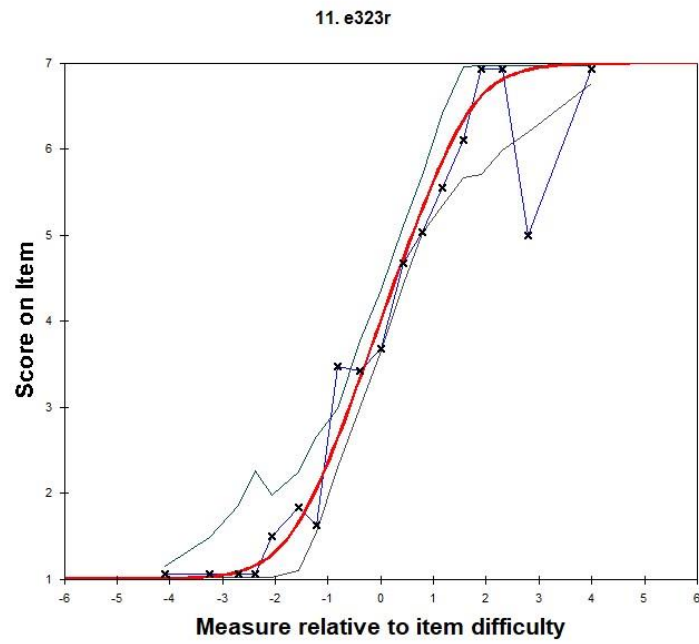
The x-axis represents the average ability of persons relative to the item score (y-axis). The solid red line is an ogive shape which shows the expected relationship between a person measure and the raw rating they would receive for that item. Areas of the curve that are steeper represent regions where the item is more discriminating The blue line represents the empirical (data-descriptive) item characteristic curve. The black "x" represents the average observed measure for an interval on the latent variable which is perceived listening effort. The grey lines represent 95% confidence intervals. The blue line is expected to approximate the red line (i.e., the empirical data fit the model) An "x" plotted outside of the confidence intervals is evidence of divergence from the Rasch-modelled expected pattern of responses"
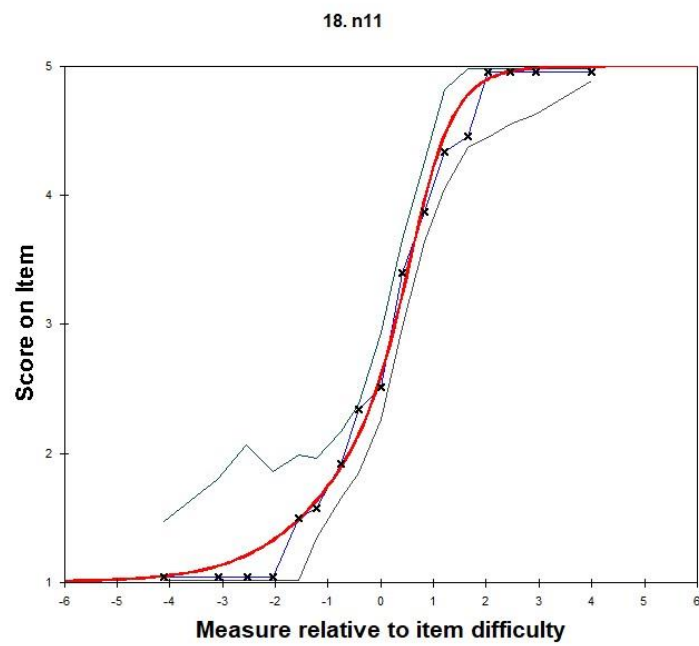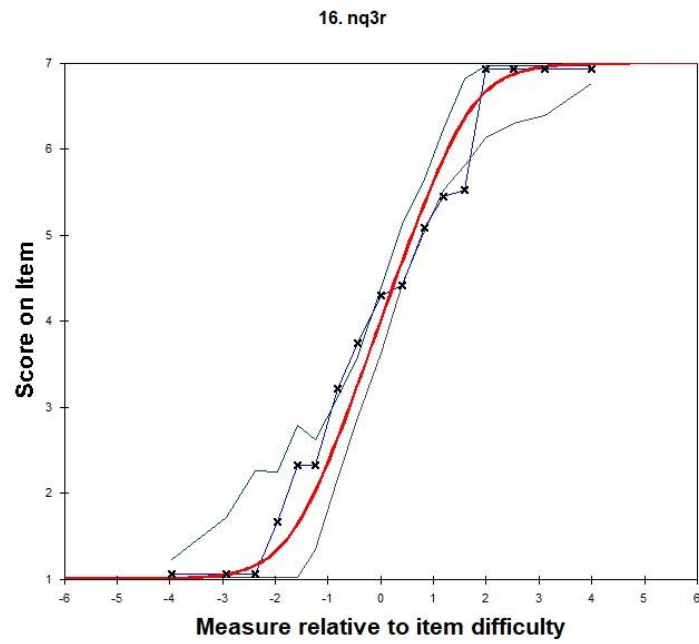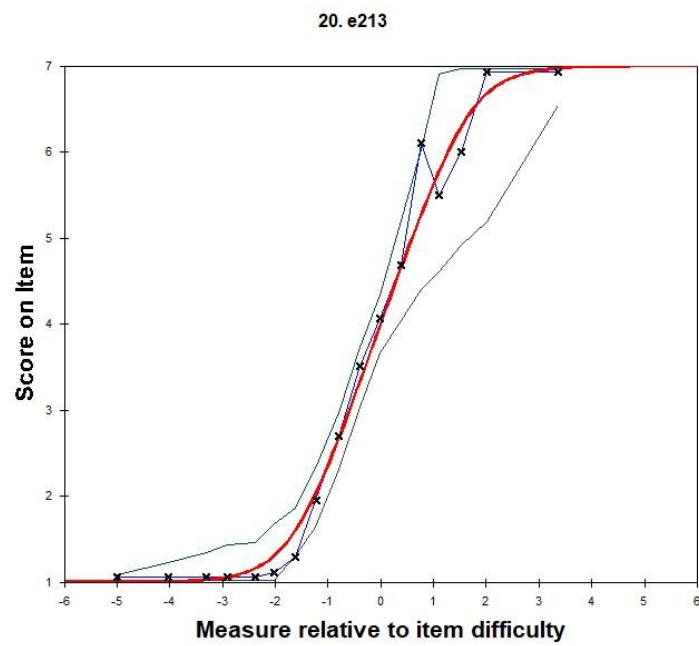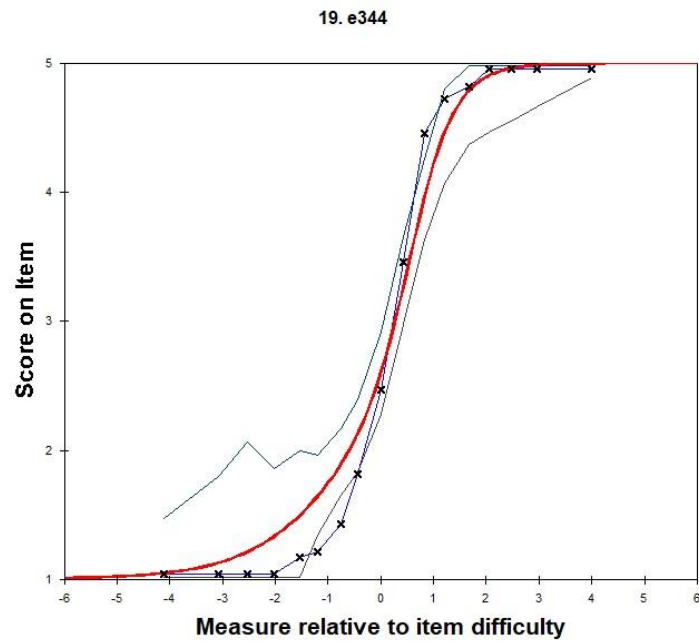


1. e398

**2. e319**



**3. nq4**

**4. e320**



**7. e35**

**9. e50**



**10. n16**

**11. e323r**



**14. e310**

**16. nq3r**



**18. n11**

**19. e344**



**20. e213**

21. e103



23. e105

**24. e214r**



**25. e118**

**26. e112**



**27. e326**

28. e363r



29. e329r

Supplemental Digital Content 4: Tables showing Winsteps generated DIF analysis (DIF) for gender and age for the LEQ-CI[21]
Table 4.1: Results of DIF analysis for gender

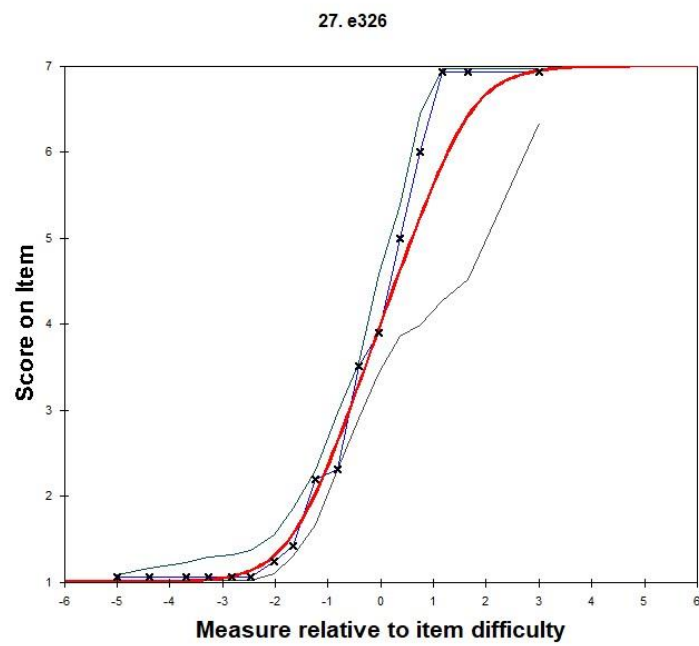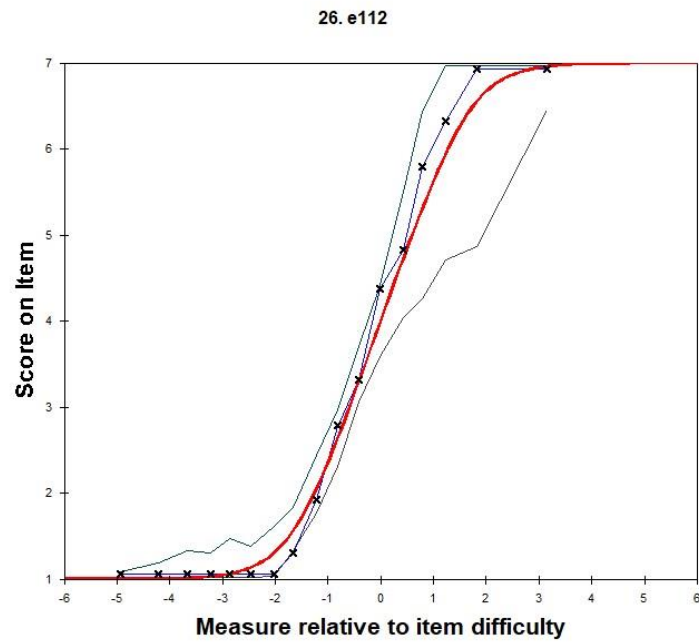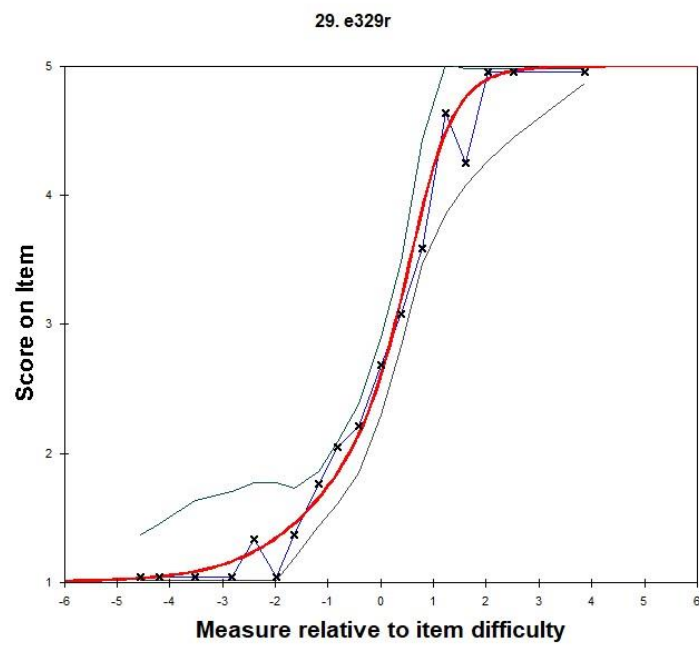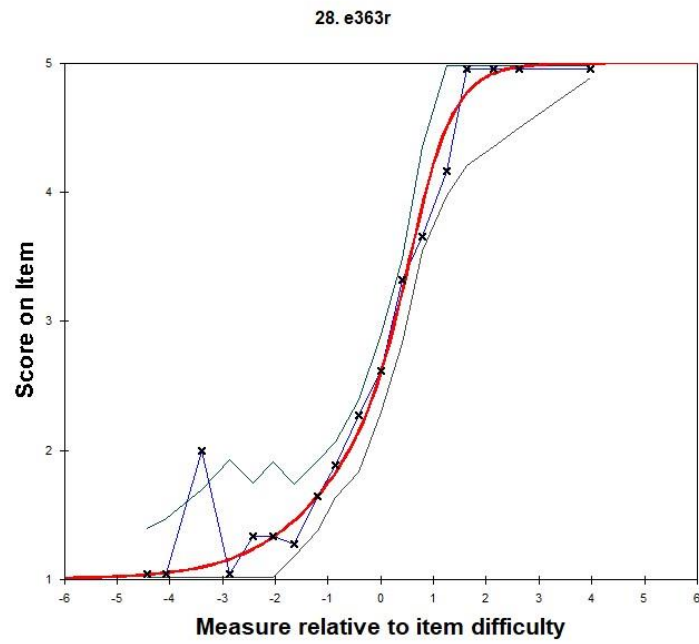| Item ID | Person Group | Observed-Expected Average | DIF Measure | DIF S.E. | Person Group | Observed-Expected Average | DIF Measure | DIF S.E. | DIF CONTRAST | Joint S.E. | Prob |
|---|---|---|---|---|---|---|---|---|---|---|---|
| e398 | M | -0.04 | -0.40 | 0.07 | F | 0.02 | -0.43 | 0.06 | 0.03 | 0.09 | 0.9408 |
| e319 | M | 0.22 | -0.42 | 0.08 | F | -0.14 | -0.09 | 0.07 | -0.33 | 0.11 | 0.0045 |
| nq4 | M | 0.02 | -0.81 | 0.08 | F | -0.01 | -0.81 | 0.07 | 0.00 | 0.11 | 0.8316 |
| e320 | M | 0.21 | -0.05 | 0.07 | F | -0.14 | 0.19 | 0.06 | -0.24 | 0.09 | 0.0063 |
| e35 | M | 0.14 | -0.34 | 0.07 | F | -0.09 | -0.19 | 0.06 | -0.15 | 0.09 | 0.0256 |
| e50 | M | 0.00 | 0.63 | 0.08 | F | 0.00 | 0.63 | 0.06 | 0.00 | 0.10 | 0.5162 |
| n16 | M | 0.03 | -0.13 | 0.07 | F | -0.02 | -0.11 | 0.06 | -0.02 | 0.09 | 0.3804 |
| e323 | M | -0.08 | -0.25 | 0.07 | F | 0.05 | -0.33 | 0.06 | 0.08 | 0.09 | 0.4454 |
| e310 | M | -0.02 | 0.60 | 0.08 | F | 0.02 | 0.60 | 0.06 | 0.00 | 0.10 | 0.6047 |
| nq3 | M | -0.04 | -0.58 | 0.07 | F | 0.02 | -0.61 | 0.06 | 0.03 | 0.09 | 0.3308 |
| n11 | M | -0.12 | -0.35 | 0.08 | F | 0.08 | -0.53 | 0.07 | 0.18 | 0.11 | 0.0408 |
| e344 | M | 0.01 | -0.47 | 0.08 | F | -0.01 | -0.47 | 0.07 | 0.00 | 0.11 | 0.8532 |
| e213 | M | 0.00 | 0.48 | 0.08 | F | 0.00 | 0.48 | 0.06 | 0.00 | 0.10 | 0.6526 |
| e103 | M | -0.05 | 1.10 | 0.08 | F | 0.04 | 1.02 | 0.07 | 0.08 | 0.11 | 0.8380 |
| e105 | M | -0.19 | 0.48 | 0.08 | F | 0.13 | 0.21 | 0.06 | 0.27 | 0.10 | 0.0268 |
| e214 | M | 0.05 | -0.29 | 0.12 | F | -0.03 | -0.15 | 0.09 | -0.14 | 0.15 | 0.8090 |
| e118 | M | -0.13 | -0.53 | 0.09 | F | 0.09 | -0.77 | 0.08 | 0.24 | 0.12 | 0.0730 |
| e112 | M | -0.04 | 0.71 | 0.08 | F | 0.03 | 0.65 | 0.06 | 0.06 | 0.10 | 0.4402 |
| e326 | M | -0.06 | 0.89 | 0.08 | F | 0.05 | 0.80 | 0.07 | 0.09 | 0.10 | 0.3312 |
| e363 | M | 0.07 | -0.20 | 0.08 | F | -0.04 | -0.09 | 0.07 | -0.10 | 0.11 | 0.3648 |
| e329 | M | 0.04 | -0.06 | 0.08 | F | -0.03 | 0.00 | 0.07 | -0.07 | 0.11 | 0.9951 |

M = males; F = females
DIF present if DIF CONSTRAST > 0.64 and p ≤ 0.05

Supplemental Digital Content 5: Raw score to 0-100 Rasch score conversion table for use with the LEQ-CI[21]

| Raw Score | Rasch Score | Raw Score | Rasch Score | Raw Score | Rasch Score |
|---|---|---|---|---|---|
| 21 | 0 | 65 | 50.81 | 109 | 68.32 |
| 22 | 12.06 | 66 | 51.19 | 110 | 68.9 |
| 23 | 18.37 | 67 | 51.56 | 111 | 69.52 |
| 24 | 21.9 | 68 | 51.92 | 112 | 70.16 |
| 25 | 24.37 | 69 | 52.29 | 113 | 70.82 |
| 26 | 26.29 | 70 | 52.65 | 114 | 71.52 |
| 27 | 27.87 | 71 | 53.01 | 115 | 72.26 |
| 28 | 29.23 | 72 | 53.37 | 116 | 73.03 |
| 29 | 30.43 | 73 | 53.72 | 117 | 73.85 |
| 30 | 31.51 | 74 | 54.07 | 118 | 74.73 |
| 31 | 32.51 | 75 | 54.43 | 119 | 75.67 |
| 32 | 33.42 | 76 | 54.78 | 120 | 76.69 |
| 33 | 34.28 | 77 | 55.13 | 121 | 77.81 |
| 34 | 35.09 | 78 | 55.48 | 122 | 79.05 |
| 35 | 35.85 | 79 | 55.83 | 123 | 80.47 |
| 36 | 36.58 | 80 | 56.18 | 124 | 82.15 |
| 37 | 37.28 | 81 | 56.53 | 125 | 84.23 |
| 38 | 37.95 | 82 | 56.89 | 109 | 68.32 |
| 39 | 38.6 | 83 | 57.24 | 110 | 68.9 |
| 40 | 39.22 | 84 | 57.59 | 111 | 69.52 |
| 41 | 39.82 | 85 | 57.95 | 112 | 70.16 |
| 42 | 40.41 | 86 | 58.31 | 113 | 70.82 |
| 43 | 40.98 | 87 | 58.67 | 114 | 71.52 |
| 44 | 41.53 | 88 | 59.04 | 115 | 72.26 |
| 45 | 42.07 | 89 | 59.4 | 116 | 73.03 |
| 46 | 42.59 | 90 | 59.78 | 117 | 73.85 |
| 47 | 43.1 | 91 | 60.15 | 118 | 74.73 |
| 48 | 43.6 | 92 | 60.53 | 119 | 75.67 |
| 49 | 44.09 | 93 | 60.92 | 120 | 76.69 |
| 50 | 44.57 | 94 | 61.31 | 121 | 77.81 |
| 51 | 45.04 | 95 | 61.7 | 122 | 79.05 |
| 52 | 45.5 | 96 | 62.11 | 123 | 80.47 |
| 53 | 45.95 | 97 | 62.52 | 124 | 82.15 |
| 54 | 46.39 | 98 | 62.94 | 125 | 84.23 |
| 55 | 46.82 | 99 | 63.37 | 126 | 87.05 |
| 56 | 47.25 | 100 | 63.8 | 127 | 91.75 |
| 57 | 47.67 | 101 | 64.25 | 128 | 100 |
| 58 | 48.08 | 102 | 64.71 | | |
| 59 | 48.48 | 103 | 65.18 | | |
| 60 | 48.89 | 104 | 65.66 | | |
| 61 | 49.28 | 105 | 66.16 | | |
| 62 | 49.67 | 106 | 66.67 | | |
| 63 | 50.06 | 107 | 67.2 | | |
| 64 | 50.44 | 108 | 67.75 | | |