1    **Development and critical evaluation of a condition-specific preference-based measure sensitive to**

2    **binaural hearing in adults: the York Binaural-related Quality of Life System**

3    A. Q. Summerfield[1], Pádraig T. Kitterick[1,2], and Adele M. Goman[1,3]

4

5    [1]Department of Psychology, University of York, York YO10 5DD, United Kingdom.

6    [2]Current address: NIHR Nottingham Biomedical Research Centre, Hearing Sciences, Division of

7    Clinical Neurosciences, School of Medicine, University of Nottingham, Nottingham NG7 2RD,

8    United Kingdom.

9    [3]Current address: Cochlear Center for Hearing and Public Health, Johns Hopkins Bloomberg School of

10    Public Health, 2024 East Monument Street, Baltimore, MD 21205, United States.

11

18

19    **Name and address for correspondence**

20    Quentin Summerfield, Department of Psychology, University of York, York YO10 5DD, UK.

21    Tel: Intl + 44 (0)1904 738825. E-mail: quentin.summerfield@york.ac.uk

22

23    **Word counts**

24    Title page (213), Abstract (499), Main Text (10,407), References and Reference Note (1,977), Figure

25    Legends, Table Legends, and Titles of Supplementary Digital Content Files (613).

26

27  **Abstract**

28  Objectives

29  The widely-used generic preference-based measures of health-related quality of life – the

30  EuroQol Descriptive System (EQ-5D) and the Health Utilities Index (HUI3) – are limited in their

31  response to technologies that improve hearing. The EQ-5D lacks construct validity for hearing, while

32  the HUI3 is restricted by a ceiling effect and by using speech reception as the only evidence of the

33  ability to hear. Consequently, neither measure consistently registers benefits from binaural hearing,

34  such as those from bilateral versus unilateral cochlear implantation. The objectives were to test

35  whether  informants value binaural hearing, to develop a condition-specific preference-based

36  measure sensitive to binaural hearing, to assess the psychometric properties of the new instrument,

37  and to determine whether it meets requirements for informing judgements of cost-effectiveness:

38  does it measure greater gains than do the generic preference-based measures, while avoiding

39  exaggerating losses, and displaying sensitivity to side effects?

40  Design

41  Three levels of function, ranging from no difficulty to great difficulty, were defined on each

42  of three dimensions where listening is easier or more successful when hearing is binaural rather than

43  monaural: perception of speech in spatially-separated noise, localization of sounds, and effort and

44  fatigue. Informants (N=203) valued the 27 combinations of levels and dimensions in a time trade-off

45  task with a 10-year time frame to provide a value of binaural-related quality of life ('binaural utility')

46  for each combination. A questionnaire was compiled to allow a respondent to report their level of

47  function on each dimension so that a value of binaural utility could be assigned to them. The

48  questionnaire and the age-standardized valuations constitute the York Binaural-related Quality-of-

49  life System (YBRQL). Adult users of unilateral implants (N=8), bilateral implants (N=11), or bimodal

50  aiding (N=9) undertook performance tests of spatial listening and completed the HUI3, EQ-5D, and

51     Speech Spatial and Qualities of Hearing (SSQ) questionnaires. They completed the YBRQL

52     questionnaire 24 and 38 months later.

53     <u>Results</u>

54     Despite long intervals between measurements, the YBRQL demonstrated desirable

55     psychometric properties: good construct validity evidenced by significant correlations with

56     performance measures and the SSQ; a greater ability than the EQ-5D or HUI3 to distinguish

57     unilateral, bimodal, and bilateral listening; and good reproducibility. The YBRQL did not exaggerate

58     losses of utility but showed no sensitivity to one potential side effect (pain/discomfort). It measured

59     a gain in utility from bilateral compared with unilateral implantation (median=.11, IQR .03 to .16)

60     that was greater than the gain measured by the EQ-5D (.00, .00 to .00) but not the HUI3 (.00, .00 to

61     .17).

62     <u>Conclusions</u>

63     The YBRQL summarizes the contribution of binaural hearing to quality of life by combining

64     the functional status of a listener with the preferences of independent informants. It would be an

65     efficient clinical outcome measure. In addition, if used alongside the EQ-5D or HUI3, it would provide

66     evidence which could beneficially modulate confidence in the cost-effectiveness of interventions.

67     Further research on its sensitivity to side effects, and on the size of the gains in utility which it

68     measures, is needed to determine whether it could stand alone to inform resource-allocation

69     decisions.

70

**Introduction**

71

72       Policy makers in healthcare require estimates of the cost-effectiveness of treatments to be

73    informed by generic preference-based measures (generic PBMs) of health-related quality of life

74    (HRQL) (e.g. Drummond et al. 2000; National Institute for Health and Care Excellence (NICE) 2013;

75    Dubois 2016). Generic PBMs define dimensions, such as pain, mobility, and anxiety, on which better

76    function corresponds to better health. The dimensions are intended to be universally relevant so

77    that generic PBMs can compare the effectiveness of interventions for widely differing conditions.

78    Respondents report their level of function on each dimension. A weighting function converts the

79    reported levels into a value of HRQL on a scale where unity corresponds to full health and zero to

80    being dead. On this scale, profoundly deaf adult candidates for cochlear implantation in the US have

81    an average value of .58 (Palmer, Niparko, & Wyatt 1999), similar to values for adults who have

82    survived heart attack (.57) or suffer chronic bronchitis (.59) (Luo et al. 2009). Following unilateral

83    implantation, adults have an average value of .78, similar to sufferers of hypertension (.73), sinusitis

84    (.76), and hay fever (.79).

85       The difference in HRQL resulting from alternative treatments can be integrated over time to

86    estimate the number of quality-adjusted life years (QALYs) gained by one treatment compared with

87    another. The additional QALYs can be combined with the additional costs to determine whether the

88    cost per QALY gained (the incremental cost-effectiveness ratio) is favourable in relation to the

89    amount that policy makers are willing for health services to pay to gain a QALY (the willingness-to-

90    pay threshold) (e.g. NICE 2013; Neumann et al. 2014; Dubois 2016).

91       This article was motivated by concerns about the validity of estimates of the gain in HRQL

92    from bilateral compared with unilateral cochlear implantation in adults. The generic PBMs which

93    have been used to estimate the gain are the 3-level version of EuroQol Descriptive System (EQ-5D-

94    3L; Brooks et al. 2003) and the Health Utilities Index Mark III (HUI3; Feeny et al. 2002). They, or a

95    third method, the time trade-off (described in Methods), were used in the five studies summarised

96    in Table 1. The estimates of the gain in HRQL vary in size between studies. More than half are small

97    and not statistically significant, despite participants reporting improvements in their listening skills

98    (Summerfield et al. 2006; Smulders et al. 2016).

99                              ----- Table 1 : Estimates of gain in HRQL -----

100          Some of the variation may stem from differences between studies in their susceptibility to

101    biases which amplify or attenuate estimates of HRQL. Retrospective studies (Bichey & Miyamoto

102    2008) are prone to selection and recall biases. Scenario analyses (Summerfield et al. 2002;

103    Kuthubutheen et al. 2015) may fail to register the impact of negative side effects. The two

104    randomized controlled trials (Summerfield et al. 2006; Smulders et al. 2016) were under-powered to

105    detect between-group differences in HRQL that were small in relation to the within-group variability.

106    Supplementary Digital Content 1 expands on this critique.

107          The small size of many of the estimates and their lack of statistical significance also reflects

108    limitations in the design of the generic PBMs. The five dimensions of the EQ-5D-3L – *Mobility, Self-*

109    *care, Usual activities, Pain/discomfort,* and *Anxiety/depression* – are minimally sensitive to sensory

110    disorders including hearing loss (Longworth et al. 2014) and to interventions which improve hearing

111    such as acoustic hearing aids (Barton et al. 2004; Grutters et al. 2007) and unilateral cochlear

112    implants (Summerfield & Barton 2019). Unsurprisingly, the EQ-5D-3L is also insensitive to the

113    difference between one implant and two.

114          In comparison, the HUI3 is consistently sensitive to hearing loss (Longworth et al. 2014) and

115    to interventions which alleviate it, including acoustic hearing aids (Barton et al. 2004; Grutters et al.

116    2007) and unilateral implants (e.g. Palmer et al. 1999; UK Cochlear Implant Study Group 2004;

117    Summerfield & Barton 2019). Sensitivity occurs because the eight dimensions ('attributes') include

118    *Hearing* (the ability to 'hear conversation') and *Speech* (the ability to be understood when speaking)*.*

119    However, the sensitivity of the HUI3 to a second implant is limited by a ceiling effect which arises

120    because the two highest levels of the *Hearing* attribute (Levels 1 and 2) require the ability to hear

121 conversation without the use of an implant/hearing aid (Horseman et al. 2003; Supplementary

122 Digital Content 2). Thus, the highest level which a user of an implant can attain is Level 3. However,

123 two thirds of users of unilateral implants place themselves at this level (Summerfield & Barton

124 2019). Therefore, only one third have sufficient headroom to report a higher level of the *Hearing*

125 attribute when using two implants rather than one.

126 Faced with a range of estimates of varying quality, systematic reviews have concluded that

127 the gain in HRQL associated with bilateral compared with unilateral implantation is sufficiently small

128 (.03, Bond et al. 2009; .035, Health Quality Ontario 2018) that there is uncertainty about the

129 likelihood that bilateral implantation is a cost-effective intervention (Bond et al. 2009; Lammers et

130 al. 2011; Health Quality Ontario 2018; Theriou et al. 2019). Analysts have reached those conclusions

131 both when they have decided on balance that bilateral implantation is a cost-effective intervention

132 for adults (Health Quality Ontario 2018) and when they have decided that it is not (Bond et al. 2009).

133 In turn, policy on the provision of bilateral implants varies within countries (Health Quality Ontario

134 2018; Boys Town National Research Hospital 2020) and between them (Vickers et al. 2016).

135 Such variation partly reflects differences between healthcare regimes in the costs of

136 provision and in willingness-to-pay thresholds, but also reflects uncertainty about the size of the gain

137 in HRQL. Herein lies a dilemma. If the true gain is of the order of .03, the likelihood is low that

138 bilateral implantation for adults achieves acceptable value for money in many healthcare

139 jurisdictions. If the gain is greater, but the generic PBMs are not equipped to detect it, then the

140 decision not to support bilateral implantation for adults risks inefficiency in the allocation of health-

141 care resources. That issue motivated the current study in which we developed and evaluated a

142 condition-specific PBM designed to be sensitive to binaural hearing in adults.

143 Condition-specific PBMs

144 The insensitivity of generic PBMs to some conditions which self-evidently reduce quality of

145 life and to some interventions which improve it has been noted in domains in addition to hearing

146 (Longworth et al. 2014). A possible solution has been the development of condition-specific PBMs

147 (e.g. Brazier et al. 2008; Yang et al. 2011; Versteegh et al. 2012; Swinburn et al. 2013). They include a

148 limited number of dimensions on which function is impaired by a condition and alleviated by an

149 effective intervention. Compared with a generic PBM that is insensitive to a condition and its

150 treatments, a condition-specific PBM would be expected to demonstrate a greater gain in HRQL

151 associated with successful treatments for the condition.

152     In practice, that advantage is in tension with weaknesses which may result in condition-

153 specific PBMs over-estimating the effectiveness of interventions (Brazier & Tsuchiya, 2010;

154 Versteegh et al. 2012). Key among them are 'exaggeration of losses' and 'insensitivity to side

155 effects'. Exaggeration of losses arises if, by focussing on a limited region of the health space,

156 differences in value within the space are amplified in the minds of informants;  as a result, greater

157 differences in value are recorded between the best and worst levels of dimensions in a condition-

158 specific PBM than would be recorded were the same dimensions valued in the context of a set of

159 generic dimensions. The second weakness, insensitivity to side effects, arises if a condition-specific

160 PBM includes dimensions chosen only to reflect the benefits, but not the disbenefits, of

161 interventions for a condition. Designers of condition-specific PBMs need to demonstrate that their

162 new instruments display the desired advantages while avoiding those weaknesses. To that end,

163 Versteegh et al. (2012) proposed criteria which condition-specific PBMs should meet before

164 informing resource-allocation decisions. In the General Discussion, we assess the condition-specific

165 PBM which we developed against those and other criteria.

166 <u>Phases of the study</u>

167     The study involved four phases. In Phase 1, we identified dimensions on which listening is

168 easier or more successful when hearing is binaural rather than monaural. We defined discrete levels

169 of function on each dimension ranging from very poor to very good. Informants used the time trade-

170 off technique (TTO; e.g. Drummond et al. 2000; York Health Economics Consortium 2016) to value

171  states of binaural hearing defined by different combinations of the levels. The TTO is a choice-based

172  method for eliciting values of states of health. It was used in many countries to value states in the

173  EuroQol Descriptive System (Brooks et al. 2003). An informant imagines that a description of a less-

174  than-perfect state applies to them and then arbitrates between two choices: either to live for a

175  specified length of time (the time frame, $y$ years) in the imperfect state, or to trade length of life for

176  quality of life and live for a shorter time ($x$ years) in full health. The informant's task is to indicate

177  what the shorter time should be for them to be indifferent between the two choices. The value of

178  health-related quality of life ('health utility') assigned by the informant to the state is calculated as

179  $x/y$. In this way, values are elicited on a scale where zero corresponds to being dead and unity to full

180  health. For example, consider that the time frame is 10 years and an informant judges that living 7.5

181  years in full health would be equivalent to living 10 years in a particular imperfect state of health.

182  The utility assigned by the informant to the imperfect state is .75 (=7.5/10). In effect, the informant

183  has equated the number of QALYs that they would live under each alternative: either 10 years at

184  75% of full health, yielding 7.5 QALYs, or 7.5 years at 100% of full health, also yielding 7.5 QALYs.

185  In our implementation of the TTO, informants imagined that descriptions of states of

186  imperfect binaural hearing applied to them. They indicated the number of years living without the

187  problems described in a state that would be equivalent to living for the time frame with the

188  problems. In this way, they assigned a value of binaural-related quality of life ('binaural utility') to

189  each state on a scale where unity corresponds to normal binaural hearing and zero to binaural

190  hearing so bad as to be equivalent to being dead.

191  There is no agreed implementation of the TTO (Attema et al. 2013), although procedures

192  used in valuing the EQ-5D-3L (Oppe et al. 2016) have been recommended (NICE 2013). Those

193  procedures include a 10-year time frame. Accordingly, we used a 10-year time frame to obtain the

194  primary dataset.

195        In Phase 2, we age-standardized the valuations. The need for age standardization arises from

196        evidence of a systematic association between age and TTO valuations. For example, when

197        informants valued health states in the EuroQol Descriptive System (Dolan & Roberts 2002), the

198        number of years traded decreased as age rose to 45 years, then increased gradually as age rose to

199        70 years and more steeply above that age. Potentially, studies would yield different average

200        valuations of the same health states if they recruited informants with different distributions of age.

201        Accordingly, we first confirmed that values of binaural utility varied with age. We then weighted the

202        data of each informant by the proportion of adults in the UK of the same age as the informant

203        before averaging the weighted data. In that way we estimated the average values of binaural utility

204        that would be obtained from a sample of informants who were representative of the distribution of

205        age in the population, rather than being tied to the particular distribution of ages in our sample.

206        In Phase 3, we compiled a questionnaire which allowed a respondent to report the level of

207        difficulty that they experienced on each dimension.

208        In Phase 4, we administered the questionnaire to adult users of cochlear implants.

209        Participants were assigned values of binaural utility according to the combination of levels they

210        adopted to describe their hearing. We employed these values to assess the potential of the

211        questionnaire in two roles: as a clinical outcome measure and as an instrument for informing

212        estimates of cost-effectiveness.

213        In summary, the aims of the study were, first, to confirm that people are willing, in principle,

214        to trade length of life to improve quality of binaural hearing; second, to develop a condition-specific

215        preference-based measure sensitive to binaural hearing in adults; third, to assess the psychometric

216        properties of the new instrument; and fourth to determine whether it meets criteria for informing

217        clinical and economic decisions in healthcare.

218    Research ethics

219        Ethical approval was granted by research ethics committees of the National Health Service

220    and the Department of Psychology of the University of York.

221

222    **Phase 1: Dimensions, levels, and values of binaural utility**

223    Methods

224    *Selection of dimensions*

225        We sought to identify dimensions with four characteristics. They should relate to differences

226    between monaural and binaural listening. They should describe fundamental (i.e.

227    physiological/psychological) and universal (i.e. occurring to the great majority of people) differences.

228    They should potentiate participation across a spectrum of activities. They should relate directly to

229    the aims of treatments and to the benefits experienced by patients. As such, they should be

230    supported by evidence of relevance from researchers, patients, and manufacturers of hearing

231    technologies. A subsidiary consideration was that informants should value every combination of

232    levels on the chosen dimensions in a single session.  That constraint limited the number of

233    dimensions and the number of levels on each dimension to three, giving 27 (=$3^3$) combinations of

234    levels and dimensions.

235        To identify dimensions, we reviewed evidence from performance tests, self-report, and the

236    guidance given by manufacturers of implants and hearing aids. Implanting or aiding both ears rather

237    than one ensures that the physiologically more responsive ear is stimulated and, potentially, gives

238    listeners access to inter-aural differences in timing and level. As a result, accuracy of localization in

239    azimuth improves – shown for bilateral compared with unilateral implantation (Kerber & Seeber

240    2012), for bimodal aiding (i.e. a unilateral implant combined with a contralateral acoustic hearing

241    aid) compared with unilateral implantation (Potts et al. 2009), and for bilateral compared with

242   unilateral fitting of acoustic hearing aids (Byrne & Noble 1998). A second consequence is that

243   accuracy of speech perception in noise improves, particularly when the sources of noise and speech

244   are spatially separated – shown for bilateral compared with unilateral implantation (Litovsky et al.

245   2006), for bimodal aiding compared with unilateral implantation (Ching et al. 2004), and for bilateral

246   compared with unilateral fittings of acoustic hearing aids (Dawes & Munro 2013).

247       Those advantages are echoed in self-reports when patients use the Speech, Spatial, and

248   Qualities of Hearing Scale (SSQ; Gatehouse & Noble 2004) to indicate how well they perform tasks of

249   speech hearing (in quiet and noise), spatial hearing (localisation of sound sources, distance and

250   movement perception), and other qualities of hearing (clarity, separation, and identification of

251   sound sources). Provision of a second implant has been associated with more positive reports in all

252   three domains (Summerfield et al. 2006; Noble et al. 2008; Härkönen et al. 2015; Smulders 2016),

253   while provision of a second hearing aid has been associated with more positive reports of speech

254   hearing and spatial hearing, and of the 'listening-effort' sub-domain of qualities of hearing (Noble &

255   Gatehouse 2006). Reductions in listening effort associated with binaural listening have also been

256   reported by users of implants (Hughes & Galvin 2013; Noble et al. 2008; Härkönen et al. 2015) and

257   demonstrated in performance tests by listeners with normal hearing (Rennies & Kidd 2018).

258       Summaries of the foregoing advantages of binaural stimulation by manufacturers of

259   implants (e.g. Advanced Bionics 2021; Cochlear 2021; Med-El 2021) and acoustic hearing aids (e.g.

260   Starkey 2021; Oticon Medical 2021) include some or all of improved abilities to localize sounds and

261   to perceive speech in noisy environments, along with enhanced clarity and greater ease of listening.

262   Those, therefore, are fundamental advantages that patients are given to expect, and often

263   experience, from bilateral fittings. Faced with the pragmatic requirement to select three dimenions,

264   we chose *Speech Perception in Noise* (SpiN), *Localization* (Loc),and *Effort and Fatigue* (E&F). We

265   judged that levels of difficulty on those dimensions could be described more simply, and would be

266     perceived to be more widely relevant, than would levels of difficulty in the fourth candidate

267     dimension, perception of sound quality.

268     *Definitions of levels*

269     We defined three levels of function on each dimension corresponding to 'No difficulty' (Level

270     1), 'Some difficulty' (Level 2), and 'Great difficulty' (Level 3). Each level was described by a brief

271     vignette (Table 2). Twenty seven scenarios were formed by combining three vignettes, one for each

272     dimension, to describe a particular combination of difficulties with binaural hearing.  Considering the

273     vignettes in the order SPiN, Loc, E&F, scenarios can be described by three numbers, referring to the

274     difficulty on each dimension, ranging from 1:1:1 to 3:3:3.

275                                          ----- Table 2 : Vignettes -----

276     Four considerations guided the construction of vignettes.  They should involve

277     straightforward language. The same wording should describe a state for an informant in a valuation

278     task and should allow a patient to report their state in a questionnaire. The wording should be

279     applicable to any condition that impairs binaural hearing and to any treatment intended to improve

280     binaural hearing. Vignettes should provide sufficient detail to define states clearly and completely.

281     The fourth consideration was addressed by describing a limitation of hearing (e.g. "You have

282     some difficulty working out where sounds are coming from."), explaining how the limitation might

283     be manifest (e.g. "You can usually tell if a sound is coming from the right- or left-hand side, but you

284     cannot be more accurate than that."), and setting out a consequence in the form of a restriction of

285     everyday activities or a hurdle to be overcome (e.g. "As a result, you are not always sure who is

286     speaking when you are in a group with several people."). In the language of the International

287     Classification of Functioning, Disability, and Health (World Health Organisation 2001), each

288     dimension relates to a *body function,* the consequences of whose impairment are illustrated by

289     restrictions on *activity* and *participation.* Supplementary Digital Content 5 argues that this strategy

290    provided complete evidence of each level of difficulty as evidenced by the statistical equivalence

291    (Lakens 2017) of the valuations of clinicians and members of the public.

292    *Valuations*

293        In each experiment, participants received a booklet containing a consent form, a

294    demographic questionnaire, instructions, and examples. The demographic questionnaire established

295    the participant's age, gender, and experience of hearing loss. Thereafter, each page contained one

296    scenario. Participants were instructed to imagine that the scenario described their own hearing and

297    to value it using a version of the TTO. Each participant valued all 27 scenarios which were presented

298    in four different randomised orders counterbalanced across participants. An experimenter was on-

299    hand to answer questions. Examples of response booklets for Experiments 1a and 1b are included in

300    Supplementary Digital Content 3 and 4.

301    *Experiment 1a*

302        Participants imagined they had 10 years to live. They decided how many years *(x)* living with

303    no hearing problems would be equivalent to living the full 10 years with the difficulties described in

304    the scenario. They indicated their choice by marking a visual-analogue scale which ranged from 0 to

305    10 with tick marks for each half year. The value of binaural utility assigned to the scenario by the

306    participant was calculated as $x/10$. The vignettes making up each scenario were presented in a fixed

307    order: SpiN, Loc, E&F.

308    *Experiment 1b*

309        In Experiment 1b, the order of presentation of the vignettes was counterbalanced among

310    participants. In addition, half of the participants made valuations while considering that their sight

311    was severely impaired while the other half considered that their sight was perfect. Here we report

312    results only for participants who considered that their sight was perfect. For them, the TTO task was

313    the same as in Experiment 1a.

314　*Participants*

315　　　　Participants (Table 3) were convenience samples of students from the University of York

316　(Students) and members of the public who were adult friends and family of students (Non-students).

317　Students divided approximately equally between those majoring in psychology and those majoring

318　in other disciplines.

319　　　　　　　　　　　----- Table 3 : Participants in Experiments 1a and 1b -----

320　*Data cleaning*

321　　　　Two of 2997 valuations were missing in Experiment 1a and 3 of 2592 in Experiment 1b.

322　Missing values were imputed deterministically based on the non-missing data such that *Imputed*

323　*value = (S x P)/G,* where *S* was the average utility for the scenario calculated from data from other

324　participants in the group, *P* was the average utility for the participant averaged over all scenarios,

325　and *G* was the grand mean utility.

326　　　　Some participants made inconsistent judgements insofar as they gave a lower value to the

327　most advantageous scenario 1:1:1 than to the least advantageous 3:3:3. In line with some studies

328　reviewed by Attema et al. (2013), inconsistent traders were excluded because they were likely to

329　have misunderstood the instructions. Participants who declined to trade when valuing any scenario

330　(zero traders) were not excluded, given that there is no reason to expect participants to consider

331　that difficulties with binaural hearing are impactful enough to justify trading length of life to alleviate

332　them. Table 3 lists the numbers of inconsistent traders, zero traders, and participants included in

333　analyses.

334　*Derived variables*

335　　　　The core data consisted of the 27 *binaural utilities* assigned by each participant. Three

336　additional measures were calculated for each participant. The average of the 27 binaural utilities

337　was calculated so that an *overall utility* could be associated with each participant as a measure of

14

338     their willingness to trade. A *mean utility* for each level of each dimension was calculated as the

339     average of the binaural utilities for the 9 scenarios in which the dimension was at a particular level. A

340     measure of the *influence* of each dimension was calculated as the difference in mean utility between

341     Level 1 and Level 3 of each dimension.

342     *Analysis*

343         Analyses were conducted with IBM SPSS for Windows v. 26.0 (2019). Effects of group,

344     dimension, and level on binaural utility were assessed in analyses of variance (ANOVAs), as were

345     effects of group and dimension on influence.  Degrees of freedom were adjusted with Huyn-Feldt

346     corrections if Mauchly's test demonstrated that the assumption of sphericity was violated.

347

348     **Results**

349     <u>Effects of Group, Dimension, and Level</u>

350     *Experiment 1a*

351         The upper row of panels in Figure 1 show how mean utility varied with *Group* and *Level* in

352     Experiment 1a. Binaural utilities were analysed in an ANOVA with the between-subjects factor *Group*

353     (Students, Non-students) and a 3x3x3 arrangement of the within-subjects factors *SpiN, Loc, E&F,*

354     each with three *Levels*. There was a significant effect of *Group* ($F_{(1,108)}$=16.617, p<.001, $\eta_p^2$=.133).

355     Students assigned lower utilities (overall utility .738, 95% confidence interval .703 to .773) than non-

356     students (overall utility .844, .806 to .882). There was also a significant effect of *Level* on mean utility

357     for each dimension: *SpiN* ($F_{(1.66,178.96)}$=100.7, p<.001, $\eta_p^2$=.482); *Loc* ($F_{(1.37,147.79)}$=135.4, p<.001, $\eta_p^2$=.556); and *E&F* ($F_{(1.37,148.01)}$=155.4, p<.001, $\eta_p^2$=.590). Mean utility declined as level varied from 1 to

359     2 and from 2 to 3 on each dimension (all p<.001). The effect of level was greater for students than

360     non-students shown by significant interactions between *Level* and *Group: SpiN* ($F_{(1.66,178.96)}$=10.363,

361    p<.001, $\eta_p^2$=.088); *Loc* ($F_{(1.37,147.79)}$=21.013, p<.001, $\eta_p^2$=.163); *E&F* ($F_{(1.37,148.01)}$=4.267, p<.001, $\eta_p^2$=.038).

362    $\eta_p^2$=.038).

363                                  ----- Figure 1 -----

364        The heights of the bars in the upper panel of Figure 2 plot the *Influence* of each dimension.

365    These measures were compared in an ANOVA with the between-subjects factor *Group* (Non-

366    students, Students) and the within-subjects factor *Dimension* (SpiN, Loc, E&F). There was a

367    significant effect of *Dimension* ($F_{(1.73,187.16)}$=22.80, p<.001, $\eta_p^2$=.174). The *Influence* of E&F (.139, .119

368    to .160) was greater than the *Influence* of Loc (.106, .090 to .123) (p<.05) which was greater than the

369    *Influence* of SpiN (.074, .062 to .087) (p<.001). There was also a significant effect of *Group*

370    ($F_{(1,108)}$=21.51, p<.001, $\eta_p^2$=.166). Students displayed a larger influence (averaged over the three

371    dimensions, .136, .119 to .153) than did non-students (.077, .059 to .096).

372                                  ----- Figure 2 -----

373    *Experiment 1b*

374        The lower row of panels in Figure 1 show how mean utility varied with *Group* and *Level* in

375    Experiment 1b*.* Binaural utilities were analysed in an ANOVA with the between-subject factor *Group*

376    (Non-students, Students) and a 3x3x3 arrangement of the within-subject factors *SpiN, Loc, E&F*, each

377    with 3 *Levels*. The effect of *Group* was close to the level of significance ($F_{(1,91)}$=3.90, p=.051, $\eta_p^2$=.041).

378    $\eta_p^2$=.041). Overall utility was .747 (.703 to .790) for students and .809 (.764 to .854) for non-students.

379    There was a significant effect of *Level* on utility for each dimension: *SpiN* ($F_{(1.47,133.96)}$=133.62,

380    p<.001, $\eta_p^2$=.595); *Loc* ($F_{(1.74,157.94)}$=121.74, p<.001, $\eta_p^2$=.572); and *E&F* ($F_{(1.41,128.54)}$=139.97, p<.001, $\eta_p^2$=.606).

381    $\eta_p^2$=.606). Mean utility declined as level varied from 1 to 2 and from 2 to 3 on each dimension (all

382    p<.001). The effect of *Level* was greater for students than non-students shown by significant

383    interactions between *Level* and *Group*: *SpiN* ($F_{(1.47,133.96)}$ =8.29, p=.001, $\eta_p^2$=.083); *Loc* ($F_{(1.74,157.94)}$=7.16,

384    p=.002, $\eta_p^2$=.073); *E&F* ($F_{(1.41,128.54)}$=4.39, p=.025, $\eta_p^2$=.046).

385    The heights of the bars in the lower panel of Figure 2 plot the *Influence* of each dimension.

386    These measures were compared in an ANOVA with the between-subjects factors *Group* (Non-

387    students, Students) and *Order* of vignettes in the scenarios (6 orders), and the within-subjects factor

388    *Dimension* (SpiN, Loc, E&F). There was a significant effect of *Group* ($F_{(1,81)}$=11.68, p=.001, $\eta_p^2$=.126);

389    students showed a larger influence (averaged over the three dimensions .129, .111 to .146) than

390    non-students (.084, .066 to .103). Influence did not vary significantly with *Order* ($F_{(5,81)}$ = 1.01) nor

391    were there any significant higher-order interactions among *Group, Sight,* and *Order.* Finally, there

392    was a significant effect of *Dimension* ($F_{(2.00, 162.00)}$=5.02, p<.01, $\eta_p^2$=.058). The influence of *E&F* (.123,

393    .104 to .141) was greater than the influence of *Loc* (.100, .084 to .115) (p<.01) and of *SpiN* (.097, .082

394    to .112) (p<.05) which did not differ significantly.

395    *Differences between Experiments*

396    Differences in overall utility between Experiments 1a and 1b were compared in an ANOVA

397    with the between-subject factors *Experiment* (1a, 1b) and *Group* (Non-students, Students). There

398    was a significant effect of *Group* ($F_{(1,199)}$=17.255, p<.001, $\eta_p^2$=.080), but not of *Experiment*

399    ($F_{(1,199)}$=.411, p=.522, $\eta_p^2$=.002).

400    Discussion

401    Three results were common to Experiments 1a and 1b. First, the majority of participants

402    were willing to trade years of life to improve quality of binaural hearing. Second, students traded

403    more years than did non-students. Third, participants traded more years to rectify problems with

404    *E&F* than problems with *SpiN* or *Loc*. In these respects, Experiments 1a and 1b replicated the results

405    of a Pilot Experiment (Supplementary Digital Content 5) which used a different implementation of

406    the TTO, thereby demonstrating that the three results are independent of the cognitive demands of

407    any particular implementation of the TTO. Finally, the absence of an effect of the order in which

408    scenarios were presented in Experiment 1b, and the absence of differences in overall utility between

409    Experiments 1a and 1b, demonstrated that the greater influence of *E&F* than *Loc* or *SpiN* in

410 Experiment 1a was not a consequence of the fixed order of presentation of the vignettes. In further

411 analyses, data from Experiments 1a and 1b were pooled.

412

413 **Phase 2: Deriving an age-standardized value set**

414 We hypothesized that a quadratic function would describe the variation in overall utility

415 with age and would accommodate the difference between students and non-students without the

416 need to invoke other explanatory variables. That hypothesis was tested by predicting overall utility

417 in Experiments 1a and 1b with a linear weighted combination of the variables *Gender* (female, male),

418 *Group* (student, non-student), and the covariates *Age* and *Age²*. Significant components of the

419 variance were explained by *Age* ($F_{(1,197)}$ = 12.152, p < .01, $\eta_p^2$ = .058) and *Age²* ($F_{(1,197)}$ = 10.959, p <

420 .01, $\eta_p^2$ = .053), but neither by *Gender* ($F_{(1,197)}$ = .456, p = .500, $\eta_p^2$ = .002), nor, critically, by *Group*

421 ($F_{(1,197)}$ = .677, p = .412, $\eta_p^2$ = .003). This result justifies combining data from students and non-

422 students after age-standardization. Further analyses of the relationship between age and overall

423 utility are reported in Supplementary Digital Content 7.

424 Table 4 illustrates the steps entailed in deriving a value set. Participants were partitioned

425 into age decades. Binaural utilities were averaged within each decade. Each average was weighted

426 by the proportion of the adult population of the UK (i.e. the population from which participants had

427 been drawn) in that decade. The weighted averages were summed. This procedure was conducted

428 for each of the 27 binaural utilities and for the overall utility.

429 ----- Table 4 : Calculation of age-standardized binaural utilities -----

430 The underlined entries in Table 5 are the 27 values that resulted, one for each scenario.

431 Values range from 0.96 to 0.69 and are well-behaved in that they vary monotonically with level

432 throughout the table.

433 ----- Table 5 : Initial and expanded value set -----

434 **Phase 3: The York Binaural-related Quality of Life System**

435      In Phase 3 we compiled a questionnaire that could elicit a person's level of function on each

436 dimension so that a value of binaural utility could be assigned to summarise the quality of their

437 binaural hearing. We named the resulting combination of questionnaire and valuations the York

438 Binaural-related Quality of Life System (YBRQL).

439 *YBRQL Questionnaire*

440      The YBRQL questionnaire is illustrated in Figure 3. The preamble invites respondents to

441 indicate their level of function on the three dimensions in the listening condition which is described

442 in bold type. Respondents choose from five levels on each dimension. Three are the levels for which

443 values of binaural utility were obtained in Phase 1. They are described by the same wording as was

444 used to elicit valuations. The other two levels are intermediary between 'No difficulty' and 'Some

445 difficulty' and between 'Some difficulty' and 'Great difficulty'. The additional levels are included to

446 allow participants to provide a granular indication of their function on each dimension.

447                                   ----- Figure 3 : YBRQL Questionnaire  -----

448 *Expanded value set*

449      To accommodate the intermediary levels of function, additional values of binaural utility

450 were calculated by linear interpolation so as to expand the 3x3x3 matrix of 27 values (underlined

451 entries in Table 5) into a 5x5x5 matrix of 125 values (all entries in Table 5). The in-between levels

452 were numbered 1.5 and 2.5 in order to preserve the numbering of the original levels.

453

454 **Phase 4: Evaluating the YBRQL (Experiment 2)**

455      We assessed three psychometric properties: *Construct validity* – does the YBRQL measure

456 difficulty with binaural listening as intended? *Discriminative ability* – is it sensitive to differences

457 between conditions, both within and between subjects? *Reproducibility* – how similar are repeated

458 measures *(Agreement)* and how consistently do they distinguish patients *(Reliability)* (discussed

459 further in Supplementary Digital Content 8). We then determined whether the YBRQL *avoided*

460 *exaggerating losses*, while measuring *greater gains* in utility than the HUI3 and EQ-5D-3L and

461 displaying sensitivity to a possible *negative side-effect.*

*Participants*

463 Participants were 28 adult users of cochlear implants, domiciled in the UK, who volunteered

464 in response to an advertisement in the newsletter of the National Cochlear Implant Users

465 Association. They were established users of cochlear implants who were willing to travel to

466 participate in two days of testing. Eight (Unilateral Group, mean age 70.6 years, SD 17.6 years, 1

467 female) used a unilateral implant only (mean duration of use of implant 10.8 years, SD 3.6 years).

468 Nine (Bimodal Group, mean age 64.3 years, SD 9.4 years, 6 female) used a unilateral implant with an

469 acoustic hearing aid stimulating the contralateral ear (mean duration of use of implant 7.3 years, SD

470 3.4 years). Eleven (Bilateral Group, mean age 64.9 years, SD 8.0 years, 5 female) used bilateral

471 cochlear implants (mean duration of use of first implant 12.4 years, SD 5.4 years; mean duration of

472 use of second implant 8.4 years, SD 4.4 years). The first language of all participants was English.

*Stage I*

474 Phase 4 involved three stages. In Stage I, participants completed a battery of performance

475 tests and self-report measures (Goman 2014) over the course of two days. On one day, they used

476 their first or only implant and completed questionnaires while thinking about themselves using only

477 one device. On the other day, the Unilateral Group was re-tested using their only implant, and the

478 Bimodal and Bilateral Groups were tested using both of their devices and completed questionnaires

479 while thinking about themselves using two devices. The assignment of number of devices to days

480 was counterbalanced as far as possible given odd numbers of participants in some groups.

481 Performance tests were conducted with an AB-York Crescent of Sound (Kitterick et al. 2011)

482 consisting of nine loudspeakers positioned at a height of 1.1m in a semicircle with a radius of 149cm.

483     Participants sat at the point equidistant from each of the loudspeakers which were located at ±90°,

484     ±60°, ±30°, ±15°, and 0° azimuth. Zero degrees was directly in front of the participant; negative

485     angles corresponded to locations to the left of straight ahead. Participants made responses on a

486     touch screen.

487         To measure the ability to localise, the phrase "Hello what's this?", spoken by an adult

488     female, was presented from one of the loudspeakers positioned, in different conditions, at (1) -60°,

489     0°, or +60°, (2) -60°, -30°, 0°, +30°,or +60°, and (3) -30°, -15°, 0°, +15°or +30°. The location was

490     chosen quasi-randomly on each trial such that each location was used equally often over the course

491     of 30 trials in each condition. The intensity of the phrase was roved from trial to trial to disrupt

492     attempts to infer the location of the source from the loudness of the stimulus at one ear.

493     Performance was scored as the percentage of trials over the three conditions on which the source

494     loudspeaker was located correctly.

495         To measure the ability to identify speech in noise, the phrase "Point to the OBJECT", spoken

496     by an adult female, was presented from the loudspeaker at 0° azimuth. OBJECT was randomly

497     selected from the set "cow", "house", "cup", "duck", "fork", "horse", "key", "tree", "plane", "plate",

498     "shoe", and "spoon". The task was to report the object named in the phrase. Pink noise was

499     presented from the loudspeaker at -90° or +90° azimuth – whichever was on the side of the

500     participant's first or only implant. This arrangement is maximally disadvantageous if hearing is

501     monaural and maximally advantageous if hearing is binaural. An adaptive procedure estimated a

502     speech-reception threshold (SRT) defined as the signal to noise ratio at which the accuracy of

503     identifying the object correctly was 70.7%. The procedure was run twice and the two estimates of

504     the SRT were averaged.

505         Participants completed the EQ-5D-3L, HUI3, and SSQ questionnaires on paper. Values of

506     health utility were obtained by using the weighting functions described by Dolan (1997) for the

507     EQ-5D-3L and Feeny et al. (2002) for the HUI3. Participants used a 0-10 visual-analogue scale to

508    respond to each of the 50 questions in the SSQ. An single score was calculated by averaging the 50

509    ratings.

510         There are two reasons why binaural utility from the YBRQL might differ from health utility

511    from the HUI3 in their relationship to other measures. First, the HUI3 is restricted by the ceiling

512    effect identified in the Introduction. Second, the HUI3 is a measure of health, not only of hearing,

513    and reflects the combined effect of poorer function on some attributes and better function on

514    others. Therefore, to compare like with like, we made an additional comparison between the YBRQL

515    and the *Hearing* attribute of the HUI3 alone, using the single-attribute utility function described by

516    Feeny et al. (2002). There is one value for each level of the *Hearing* attribute. It was obtained from

517    informants who valued each level while all other attributes were set to their highest level.  The

518    resulting values were then scaled to range from 1.00 (Level 1) to 0.00 (Level 6). We refer to these

519    values as measures of *hearing* utility, whereas we refer to the multi-dimensional values from the

520    HUI3 and EQ-5D-3L as measures of *health* utility.

521    *Stage II*

522         Twenty four months after Stage I, the YBRQL questionnaire was mailed to participants. The

523    Unilateral Group completed the questionnaire once. The Bimodal and Bilateral Groups completed

524    the questionnaire twice, first considering their functional abilities when using their first or only

525    implant (Monaural Condition), and second when using both of their devices (Binaural Condition). For

526    the monaural condition, the preamble instructed participants to "indicate which statement best

527    describes your own hearing when using your cochlear implant" (Unilateral Group), "when using your

528    cochlear implant on its own" (Bimodal Group) or "when using your first cochlear implant on its own"

529    (Bilateral Group). For the binaural condition, the preamble instructed participants to "indicate which

530    statement best describes your own hearing when using your cochlear implant together with your

531    hearing aid" (Bimodal Group) or "when using your two cochlear implants together" (Bilateral Group).

532    A value of binaural utility from Table 5 was assigned to the participant according to the levels of

533    function reported.

534    *Stage III*

535    Fourteen months after Stage II, the YBRQL questionnaire was mailed to participants who

536    completed it while considering their functional abilities with their usual configuration of devices. A

537    value of binaural utility from Table 5 was assigned to the participant according to the levels of

538    function reported.

539    *Interim summary: Outcome measures*

540    In summary, we derived values of binaural utility from the YBRQL for comparison with two

541    sets of measures. The first set consisted of the SSQ score and the performance measures of

542    localization and speech perception in noise; we regarded these as *gold-standard* measures of

543    binaural hearing. The second set consisted of hearing utility from the HUI3 and health utility from

544    the HUI3 and EQ-5D-3L; we refer to these as *reference* measures. In addition, we calculated the

545    *binaural advantage* for each measure as the difference between values obtained with one device

546    and two devices by the Bimodal and Bilateral Groups.

547    <u>Data Analysis</u>

548    Analyses were conducted with IBM SPSS for Windows v. 26.0 (2019). If all measures in a set

549    of related analyses distributed normally, parametric tests are reported. If one or more measures

550    departed from normality, evidenced by significant Shapiro-Wilk tests, non-parametric tests are

551    reported. Those reports are supplemented, if informative, by results of parametric tests for

552    measures that distributed normally. Where there was no straightforward non-parametric alternative

553    to a parametric test, such as the analysis of covariance, the results of the parametric test are

554    reported. However, conclusions are drawn only if significant results are corroborated by confidence

555    intervals estimated by bootstrapping (3,000 samples per analysis; bias corrected and accelerated).

*Construct validity*

556

557    Construct validity was assessed by calculating correlation coefficients between binaural

558    utility from the YBRQL and the three gold-standard measures. First, correlations were calculated

559    with data obtained when participants thought about, and performed with, their usual configuration

560    of devices. Second, correlations were calculated between the measures of binaural advantage. Third,

561    for comparison with the YBRQL, corresponding correlations were calculated between the gold-

562    standard measures and the reference measures. Correlation coefficients were classified as small (.1

563    to <.3), moderate (.3 to <.5), large (.5 to <.7), or very large (≥.7). Non-parametric coefficients of at

564    least moderate size, and parametric coefficients of at least large size, were taken as evidence of

565    construct validity.

*Discriminative ability*

566

567    The ability of the YBRQL to discriminate between conditions was assessed by testing two

568    hypotheses. The first was that binaural utility from the YBRQL was higher when members of the

569    Bimodal and Bilateral groups used two devices rather than one. Effect sizes were classified as small

570    (≥.2 to <.5), medium (.5 to <.8), or large (≥.8) (Cohen, 1988). The second hypothesis was that the

571    binaural advantage measured by the YBRQL was greater for the Bilateral group than the Bimodal

572    group. This hypothesis was tested in analyses of covariance (ANCOVAs) with the score with the first

573    or only implant as a covariate to control for differences among participants at baseline (Vickers &

574    Altman 2001). Corresponding analyses were conducted with each of the reference measures.

*Reproducibility*

575

576    Reproducibility was assessed for values of binaural utility from the YBRQL. *Agreement* was

577    quantified as the standard error of measurement (SEm), calculated as the square root of the sum of

578    the variance due to conditions (test and retest) and the residual variance due to the interaction

579    between conditions and participants (de Vet et al. 2006). *Reliability* was quantified as the intra-class

580    correlation coefficient (ICC) obtained when participants considered their functional abilities with

24

581      their usual configuration of devices in Stages II and III. The ICC was calculated with a 2-way mixed-

582      effects model assessing absolute agreement for single measures (McGraw & Wong 1996). Following

583      Koo and Li (2016), values of the ICC were classified as poor (<.5), moderate (.5 to <.75), good (.75 to

584      <.9), or excellent (≥.9). An ICC of at least good quality was required as evidence of reliability.

585      *Losses of utility*

586      To test whether the YBRQL exaggerated losses of utility, we compared the YBRQL with the

587      HUI3 in terms of the loss relative to normal hearing when participants listened with their first or only

588      implant. For the YBRQL, the loss of utility was calculated as the difference between the maximum

589      attainable value, .96 (Table 5), and the observed value. For the HUI3, the measure was the

590      difference between the health utility when the hearing dimension was set to its highest level while

591      all other dimensions were at their observed levels, and the health utility when all dimensions were

592      at their observed levels.

593      *Gains in utility*

594      We tested whether the YBRQL measured larger gains than the HUI3 and EQ-5D-3L by

595      determining whether the binaural advantages shown by the Bimodal and Bilateral Groups were

596      significantly larger when measured as changes in binaural utility than as changes in health utility.

597      *Sensitivity to side effects*

598      While serious and long-lasting complications associated with cochlear implantation are rare

599      (e.g. Figure 3 in UKCISG 2004) once clinical programs are experienced (e.g. Cohen et al. 1993),

600      increases on the pain/discomfort dimensions of generic PBMs have been reported in the short term

601      (Summerfield & Barton 2019), possibly arising from the surgical wound or from the receiver-

602      stimulator irritating the scalp. Also, while tinnitus is generally reduced by implantation, it may be

603      exacerbated in some patients and induced in others (Ramakers et al. 2015) and is associated with

604      changes on the Anxiety/Depression dimension of the EQ-5D (Summerfield & Barton 2019). These

605    effects, although often mild and sub-clinical, have the potential to negate gains in utility associated

606    with improved binaural hearing. The YBRQL might be insensitive to these side-effects. Alternatively,

607    they might be associated with increased effort and fatigue.

608         The data of Experiment 2 permitted only a limited assessment of sensitivity to side effects in

609    the form of a test of the sensitivity of the YBRQL to variation among participants in pain/discomfort

610    whether related to implantation or not. We conducted an ANCOVA with binaural utility as the

611    dependent variable, group as a fixed factor, and the levels of the *Hearing* and *Pain* dimensions of the

612    HUI3 as covariates. The critical test was whether both covariates were significantly related to the

613    dependent variable. For comparison, we applied the same test to health utilities estimated with the

614    EQ-5D-3L.

615    Results

616    *Response rate*

617         Of 31 participants who took part in Stage 1, 28 returned the YBRQL questionnaire in Stage 2,

618    and, of them, 25 returned it in Stage 3. Non-respondents had either died or were unwell.

619    *Construct validity*

620         The scatter plots in Figure 4 show the relationship between the SSQ score and measures of

621    utility. The relationship with binaural utility from the YBRQL (Figure 4A) is characterised by a

622    moderate-to-large non-parametric correlation (Kendall's $\tau$ = .501, N=28, p<.001) and by a large

623    parametric correlation (Pearson's r = .680, N=28, p<.001). The relationship with hearing utility from

624    the HUI3 (Figure 4B) is characterised by a moderate non-parametric correlation ($\tau$ = .463, N=28,

625    p=.003). The relationship with health utility from the HUI3 (Figure 4C) is also characterised by a

626    moderate non-parametric correlation ($\tau$ = .386, N=28, p<.01). No relationship was found with health

627    utility from the EQ-5D-3L ($\tau$ = .215, N=28, p=.146) (Figure 4D).

628       Panels E, F, G, and H show the corresponding relationships among the measures of binaural

629       advantage. Binaural advantage from the YBRQL was significantly associated with binaural advantage

630       from the SSQ ($\tau$ = .492, N=20, p=.003; r = .752, N=20, p<.001). Binaural advantage from none of the

631       three reference measures was significantly associated with binaural advantage from the SSQ.

632       ----- Figure 4 : Scatterplots -----

633       Correlations between the gold-standard performance measures and measures of utility are

634       listed in Table 6.  Higher values of binaural utility from the YBRQL were associated with better

635       speech perception (i.e. lower SRTs) and more accurate localisation. Higher values of health utility

636       from the HUI3 were associated with better speech perception, but not with accuracy of localisation.

637       There was no association between health utility from the EQ-5D-3L and either localisation or speech

638       perception.

639       ----- Table 6 : Correlations -----

640       *Discriminative ability*

641       Figure 5 contains box plots of the SSQ score and of measures of utility for listening with one

642       device (open bars) and two devices (shaded bars) by the Bilateral and Bimodal Groups. The panel

643       beneath the plots includes the results of Wilcoxon Signed Ranks tests which compared scores with

644       one and two devices. The SSQ score and binaural utility from the YBRQL were significantly higher

645       with two devices than one, while there were no significant differences in hearing utility from the

646       HUI3 or in health utility from the HUI3 and EQ-5D-3L.

647       Effect sizes were calculated by dividing the *z-score* by the square root of the number of

648       observations (Field 2009) and are tabulated as values of *R* at the bottom of Figure 5. Medium effect

649       sizes were shown by the SSQ score and the measure of binaural utility, while either small (HUI3) or

650       negligible (EQ-5D-3L) effect sizes were shown by the other measures of utility.

651       ----- Figure 5 : Box plots -----

652     ANCOVAs tested the hypothesis that binaural advantage was greater for the Bilateral Group

653     than the Bimodal Group. *Group* (Bilateral, Bimodal) was a fixed factor. *Score with the first or only*

654     *implant* was a covariate. The effect of *Group* was significant for the SSQ score ($F_{(1,17)}$=8.933, p=.008,

655     $\eta_p^2$=.344), for binaural utility from the YBRQL ($F_{(1,17)}$=10.742, p=.004, $\eta_p^2$=.387), but not for hearing

656     utility from the HUI3 ($F_{(1,17)}$=.224, p=.642), nor for health utility from the HUI3 ($F_{(1,17)}$=1.161, p=.296)

657     or the EQ-5D-3L ($F_{(1,17)}$=.210, p=.653). This pattern of significance was corroborated by bootstrapped

658     estimates of the 95% confidence intervals of the binaural advantages (taken as the estimated

659     marginal means) and their difference (taken as the parameter for *Group* in the general linear model)

660     (Table 7).

661                          ----- Table 7 : Discriminative ability -----

662     The box plots in Figure 6 show the extent to which levels of function on the individual

663     dimensions of the YBRQL differed between monaural and binaural conditions. The Localization

664     dimension showed a significant improvement for the Bimodal Group (Wilcoxon Signed Ranks Test,

665     z=2.271, p<.05), while all three dimensions showed significant improvements for the Bilateral Group

666     (SpiN, z=2.701, p<.01; Loc, z=2.579, p<.01; E&F, z=2.565, p<.05).

667                          ----- Figure 6 : Contributions of individual dimensions ----

668

669     *Reproducibility*

670     Mean values of average utility from the YBRQL differed minimally between test, .834 (SD =

671     .068), and retest, .835 (.057). The SEm was .03. The ICC (with 95% confidence interval) was .818

672     (.628 to .916).

673     *Losses of utility*

674     The upper part of Table 8 lists median values of the losses of utility attributable to impaired

675     hearing from the YBRQL and HUI3 when participants used their first or only implant.  The losses did

676    not differ significantly (Wilcoxon Signed Ranks Test: N=28, z=1.586, exact p=.115, 2-tailed). There is

677    no evidence, therefore, that the YBRQL exaggerated losses in comparison with the HUI3.

678                                     ----- Table 8 : Losses and gains in utility -----

679    *Gains in utility*

680         The lower part of Table 8 lists median values of the binaural advantages recorded by the

681    YBRQL, HUI3, and EQ-5D-3L. For the Bimodal Group, the advantage measured with the YBRQL was

682    not greater than the advantage measured with the EQ-5D-3L (Wilcoxon Signed Ranks Test: z=.169,

683    N=9, exact p=.469, 1-tailed) or the HUI3 (z=.280, exact p=.422, 1-tailed). For the Bilateral Group, the

684    advantage measured with the YBRQL was greater than the advantage measured with the EQ-5D-3L

685    (N=11, z=2.312, exact p=.009, 1-tailed), and fell short of significance in comparison with the

686    advantage measured with the HUI3 (z=1.689, exact p=.051, 1-tailed).

687    *Sensitivity to side effects*

688         Binaural utility from the YBRQL when participants used their normal configuration of devices

689    was analysed in an ANCOVA with the fixed factor *Group* (Unilateral, Bimodal, Bilateral) and the

690    covariates *Level of HUI3 Hearing* (Hearing) and *Level of HUI3 Pain* (Pain). There were significant

691    effects of *Group* ($F_{2,23}$=7.830, p<.01, $\eta_p^2$=.405) and *Hearing* ($F_{1,23}$=9.111, p<.01, $\eta_p^2$=.284), but not of

692    *Pain* ($F_{1,23}$=.154, n.s., $\eta_p^2$=.007). In comparison, an analogous analysis of health utility from the

693    EQ-5D-3L showed the reciprocal relationship: a significant effect of *Pain* ($F_{1,23}$=13.154, p<.01,

694    $\eta_p^2$=.364), but not of *Group* ($F_{2,23}$=1.493, n.s., $\eta_p^2$=.115) or *Hearing* ($F_{1,23}$=.923, n.s., $\eta_p^2$=.039). Thus,

695    the measure of binaural utility from the YBRQL distinguishes patients according to differences in

696    their ability to hear but not according to differences in the pain/discomfort which they report. The

697    fact that the measure of health utility from the EQ-5D-3L is sensitive to differences in

698    pain/discomfort shows that the insensitivity of the YBRQL was a true limitation rather than the result

699    of negligible variation in pain/discomfort. Further analyses assessing the sensitivity of the YBRQL to

700    pain/discomfort are reported in Supplementary Digital Content 2.

701    <u>Discussion</u>

702    *Psychometric properties*

703    Binaural utility measured with the YBRQL displayed good construct validity (Figure 4; Table

704    6) and discriminative ability (Figure 5; Table 7) with all three dimensions contributing (Figure 6). In

705    these respects, it performed similarly to the SSQ, better than the HUI3, and much better than the

706    EQ-5D-3L.  It also demonstrated good reproducibility over a 14-month interval between test and

707    retest. The value of the SEm, .03, measuring agreement, compares favourably with values for the

708    EQ-5D and HUI3 (Palta et al. 2011) which were close to .10 (at an interval of 6 months by 250

709    patients following cataract surgery). The value of the ICC (.82, .63 to .92), measuring reliability, is

710    similar to values reported for the generic PBMs (EQ-5D-3L at an interval of 3 months by 224 patients

711    with rheumatoid arthritis: .82, .74 to .88; Macran, 2003; HUI3 at an interval of 1 month by 506

712    members of the general public: .77; Boyle et al. 1995; and at an interval of 3 months by 141 patients

713    with hip fracture: .77; Jones et al. 2005).

714    *Other statistical properties*

715    Binaural utility avoided floor and ceiling effects and displayed acceptable effect sizes when

716    contrasting binaural with monaural hearing (Figures 4A and 5). With a finely graded set of response

717    options to choose from, 17 of the 20 members of the Bimodal and Bilateral Groups reported an

718    improvement in binaural utility when using two devices rather than one. In contrast, hearing utility

719    measured with the HUI3 showed the ceiling effect anticipated in the Introduction (Figures 4B and 5);

720    seventeen of 28 participants reported that they functioned at ceiling (Level 3) when using their first

721    or only implant and only 6 of the 20 members of the Bimodal and Bilateral Groups reported an

722    improvement when using two devices rather than one. The same limitations restricted the sensitivity

723    of health utility from the HUI3 to binaural hearing (Figures 4C and 5), while the sensitivity of health

724    utility from the EQ-5D-3L was restricted because more than half of the participants reported levels

725    of function corresponding to full health when using one device (Figures 4D and 5).

726

**727    General Discussion**

728         There are four respects in which it is relevant to assess the YBRQL: Its potential as a clinical

729    outcome measure; its legitimacy for, and its practical use in, informing resource-allocation decisions;

730    limitations that arose in its development; and desirable further developments.

731    <u>Potential as a clinical outcome measure</u>

732         The levels of function reported by a participant with the YBRQL questionnaire describe a

733    profile across three fundamental dimensions of binaural hearing: speech perception in noise,

734    localization, and the effort and fatigue entailed in listening. The value of binaural utility that can be

735    derived from the profile provides a summary index of binaural function. The index ranges in value

736    from .69 to .96. The lower value corresponds to great difficulties on all three dimensions; the higher

737    value corresponds to normal performance on each dimension. The index results from the mental

738    integration of the impact of the difficulties by independent informants. The index displays good

739    psychometric properties of construct validity, discriminative ability, and reproducibility.  Considered

740    overall, the YBRQL provides a disciplined summary of one way of estimating the contribution of

741    binaural hearing to quality of life. It would be an efficient and sensitive clinical outcome measure in

742    studies that require a brief questionnaire that provides a quantitative summary of binaural function

743    that correlates robustly with performance measures and the overall score from the longer SSQ.

744     <u>Legitimacy for informing resource-allocation decisions</u>

745          Legitimacy hinges on whether the YBRQL meets the criteria set out by Versteegh et al.

746     (2012), including demonstrating larger gains in utility than are shown by generic PBMS, while

747     avoiding exaggeration of losses and displaying sensitivity to side effects.

748     *Assessment against criteria*

749          The first criterion proposed by Versteegh et al. (2012) is that *empirical evidence disproves*

750     *the sensitivity of existing generic instruments.* Previous studies (Table 1) and Experiment 2 (Tables 6

751     and 7; Figures 4 and 5) show that the EQ-5D-3L is largely insensitive to binaural hearing while the

752     HUI3 displays restricted sensitivity. The first criterion is partially met, therefore. The second criterion

753     is that *empirical evidence proves the superiority of the condition-specific instrument from which the*

754     *condition-specific PBM was derived.* The YBRQL did not draw questions from a pre-existing

755     instrument. However, its dimensions align approximately with the three main sections of the SSQ.

756     The SSQ is sensitive to binaural hearing in conditions where generic PBMs are either insensitive

757     (Summerfield et al. 2006; Smulders et al. 2016) or less sensitive (Figure 5; Table 8). To that extent,

758     the second criterion is met.

759          The third criterion is that the *derived condition-specific PBM is shown to be superior to the*

760     *existing generic PBM(s), not just in terms of statistical sensitivity, but also in terms of absolute*

761     *differences* (Reference Note 1). The YBRQL meets the 'statistical sensitivity' part of the criterion. It

762     displayed larger effect sizes than the HUI3 or EQ-5D-3L when comparing monaural with binaural

763     hearing (Figure 5). It displayed a consistent pattern of significant differences when monaural and

764     binaural hearing were contrasted both within and between groups, whereas the HUI3 displayed an

765     inconsistent pattern and the EQ-5D-3L showed no significant differences (Table 7). However, the

766     YBRQL failed the 'absolute differences' part of the criterion, in that the binaural advantages

767    measured by the YBRQL for the Bimodal and Bilateral groups were not significantly larger than the

768    advantages measured by the HUI3.

769          This last result is important, given that the motivation for developing the YBRQL was concern

770    about the lack of sensitivity of the HUI3 to bilateral implantation. Note, however, that the difference

771    between the binaural advantage for the Bilateral Group measured with the YBRQL (median=.107,

772    IQR .035 to .162) and the HUI3 (.000, .000 to .170) fell only marginally short of significance (p=.051

773    when assessed with a 1-tailed Wilcoxon Signed Ranks test). Further analyses (Supplementary Digital

774    Content 2) confirm that the significance of the difference fell in the region of uncertainty and that a

775    more powerful study is required to resolve the issue.

776    *Exaggeration of losses*

777          The aim of developing the YBRQL was to avoid limitations in the sensitivity of the HUI3 to

778    interventions which improve binaural hearing, rather than to increase sensitivity to impaired

779    hearing. Appropriately, therefore, the YBRQL did not exaggerate the size of losses of utility due to

780    impaired hearing when compared with the HUI3.

781    *Sensitivity to side effects*

782          Experiment 2 demonstrated that the YBRQL is sensitive to differences in the configuration of

783    implants and acoustic hearing aids. Accordingly, it would be expected also to be sensitive to

784    complications, such as the malfunctioning or sub-optimal fitting of an implant or hearing aid, that

785    reduce binaural function. However, it is not clear whether the YBRQL is also sensitive to other side-

786    effects, such as changes in pain/discomfort (Summerfield & Barton 2019) and tinnitus (Ramakers et

787    al. 2015), whose impact may be orthogonal to binaural function. Experiment 2 showed that the

788    YBRQL is insensitive to variation among participants in pain/discomfort, whether related to

789    implantation or not. That may not matter, given that pain is reported to be an unusual long-term

790    consequence of cochlear implantation (Celerier et al. 2017) and at least one study has reported no

791    worsening of pain, or any other attribute of the HUI3, associated with bilateral implantation (Table 4

792    in Bichey & Miyamoto 2008). On the other hand, two members of the Bilateral Group in Experiment

793    2 reported greater pain/discomfort when using two implants compared with one (Supplementary

794    Digital Content 2), and it would be unwise to assume that any intervention which entails surgery is

795    immune to negative side effects. Further research is needed to establish the sensitivity, or

796    otherwise, of the YBRQL to side effects.

797    *Role of the YBRQL*

798         Where do these assessments leave the YBRQL? One option would be to restrict its role to

799    providing a measure of clinical effectiveness, while the HUI3 informs analyses of cost-effectiveness.

800    However, that strategy is undermined by the evidence in Table 1 which motivated the development

801    of the YBRQL, by the intrinsic limitations in the sensitivity of the HUI3 to binaural hearing identified

802    in the Introduction and in Experiment 2, and by further limitations discussed in Supplementary

803    Digital Content 2.

804         The second option would be to employ the YBRQL to supplement evidence provided by the

805    EQ-5D or HUI3. Consider that three willingness-to-pay thresholds are relevant to resource-allocation

806    decisions. The first is a threshold below which an intervention is unlikely to be rejected on grounds

807    of cost-ineffectiveness. The second is a higher threshold above which special reasons are required if

808    an intervention is to be accepted. Implicitly, there is a third, yet higher, threshold above which

809    interventions are never accepted. The lower two thresholds may be set explicitly (e.g. Rawlins &

810    Culyer, 2004; NICE, 2013) or inferred (Neumann et al. 2014; Dubois 2016; Cameron et al. 2018). The

811    third threshold is necessarily inferred (e.g. Dakin et al. 2015). The range between the first and third

812    thresholds defines a region of uncertainty. Figure 7 is a matrix of decisions that could be reached by

813    using the YBRQL to estimate the cost-effectiveness of binaural-related gains in quality of life

814    alongside the EQ-5D or HUI3 to estimate the cost-effectiveness of health-related gains in quality of

815    life. The cells correspond to the nine possible combinations of two incremental cost-effectiveness

816    ratios each of which may be more favourable than the first threshold, or in the region of uncertainty,

817     or above the third threshold. Of particular relevance is Cell 3, where evidence from the YBRQL might

818     prompt further investigation before accepting evidence of cost-effectiveness from the generic PBM,

819     and Cell 4, where evidence from the YBRQL might prompt the judgement that an intervention was

820     cost-effective despite uncertainty in the evidence from the generic PBM.

821                              ----- Figure 7 : Decision matrix -----

822     <u>Limitations</u>

823        Experiments 1 and 2 have limitations. First, although the choice of dimensions for the YBRQL

824     was principled, the inclusion of additional dimensions informed by the ICF (WHO 2001) or relevant

825     questionnaires (e.g. Hinderink et al. 2000) might result in greater sensitivity to improvements in

826     binaural hearing. Second, the informants who contributed to the valuation of the YBRQL in

827     Experiment 1 were a convenience sample rather than a population-representative sample and

828     included relatively few older people. Third, the evaluation of the YBRQL in Experiment 2 was limited

829     to interventions involving cochlear implants. Fourth, the value set for the YBRQL was obtained from

830     informants domiciled in the UK. Valuations of health states differ between countries (e.g. Bernert et

831     al. 2009; Gerlinger et al. 2019), at least partly because different populations place different values on

832     the same aspects of health depending on the context in which healthcare is provided. For those

833     reasons, policy-makers prefer decisions to be informed by preference data from their own

834     population (e.g. NICE 2013). Researchers outside the UK considering using the YBRQL in partnership

835     with a generic PBM to inform resource-allocation decisions might choose to generate a value set

836     from a sample of their own population.

837        The fifth limitation is that, while informing analyses of the psychometric properties of the

838     YBRQL, Experiment 2 was under-powered to detect differences between gains in utility measured

839     with different PBMs, given the wide dispersion of utilities measured with the HUI3 and EQ-5D-3L

840     (Figure 5), and was not designed to examine the impact of side effects. It would be desirable to

841     recruit a larger sample of users of unilateral implants wo received a second implant, and who were

842   tested prospectively, to address those issues. If the resolution was not satisfactory, it would be

843   desirable to integrate the YBRQL with a generic PBM by adding dimensions to the generic PBM

844   assessing *Speech-perception in noise* and *Localisation*. 'Bolting on' additional dimensions has been

845   advocated where the EQ-5D is insensitive to a condition and to treatments that alleviate it (e.g.

846   Krabbe et al. 1999; Yang et al. 2015; Finch et al. 2017). Sensitivity to a condition is enhanced, while

847   avoiding exaggeration of losses and preserving the capacity to register side-effects. Although that

848   approach did not yield systematic valuations when questions asking directly about the ability to

849   'hear' were bolted on (Yang et al. 2015), the present results demonstrate that informants are willing

850   to trade, and produce systematic results, when questions more fully exemplify the aspects of

851   listening which are at stake.

852   Conclusion

853       The York Binaural-related Quality of Life System provides a measure of the contribution of

854   binaural hearing to the quality of life of adult listeners. It would be an efficient clinical outcome

855   measure in the evaluation of interventions intended to improve binaural hearing. At its present

856   stage of development, it should not be used on its own to inform resource-allocation decisions.

857   However, if administered in conjunction with a generic PBM, it could beneficially modulate

858   confidence in the cost-effectiveness of interventions.  To that end, Supplementary Digital Content 9

859   contains the YBRQL questionnaire as an editable document which can be tailored to a particular

860   application.

861

866

867 **Author contributions**

868     AQS, PTK, and AMG designed the experiments. AQS and AMG supervised data-gathering for

869 Experiment 1. AMG and PTK gathered the data for Experiment 2. AQS, PTK, and AMG analysed the

870 data. AQS and AMG wrote the paper.

871

872 **References**

873 Attema, A.E., Edelaar-Peeters, Y., Versteegh, M.M., et al. (2013). Time trade-off: one methodology,

874     different methods. *Eur J Health Econ, 14(Suppl 1)*, S53-S64.

875 Barton, G.R., Bankart, J., Davis, A.C., et al. (2004). Comparing utility scores before and after hearing-

876     aid provision. *Applied Health Economics and Health Policy, 3*, 103-105.

877 Bernert, S., Fernández, A., Haro, J.M., et al. (2009). Comparison of different valuation methods for

878     population health status measured by the EQ-5D in three European countries. *Value in Health,*

879     *12,* 750-758.

880 Bichey, B.G., Miyamoto, R.T. (2008). Outcomes in bilateral cochlear implantation. *Otolaryngology—*

881     *Head and Neck Surgery, 138,* 655-661.

882 Bond, M., Mealing, S., Anderson, R., et al. (2009). The effectiveness and cost-effectiveness of

883     cochlear implants for severe to profound deafness in children and adults: a systematic review

884     and economic model. *Health Technol Assess, 13,* 1-330.

885 Boyle, M. H., Furlong, W., Feeny, D., et al. (1995). Reliability of the Health Utilities Index-Mark III

886     used in the 1991 Cycle 6 Canadian General Social Survey Health Questionnaire. *Quality of Life*

887     *Research 4,* 249-257.

888  Boys Town National Research Hospital (2020). Statement on insurance coverage of bilateral cochlear

889  implantation. Accessed at https://www.boystownhospital.org/knowledge-center/bilateral-

890  cochlear-implants on 8th June 2020.

891  Brazier, J., Czoski-Murray, C., Roberts, J., et al. (2008). Estimation of a preference-based index from a

892  condition-specific measure: the King's Health Questionnaire. *Med Decis Making, 28,* 113-126.

893  Brazier, J., Tsuchiya, A. (2010). Preference-based condition-specific measures of health: what

894  happens to cross programme comparability? *Health Economics, 19,* 125-129.

895  Brooks, R., Rabin, R., de Charro, F. (2003). *The Measurement and Valuation of Health Status using*

896  *EQ-5D: A European Perspective.* Dordrecht: Kluwer Academic Publishers.

897  Byrne, D., Noble, W. (1998). Optimizing sound localization with hearing aids. *Trends in Amplification,*

898  *3, 51-73.*

899  Cameron D., Ubels, J., Norstrom, F. (2018). On what basis are medical cost-effectiveness thresholds

900  set? Clashing opinions and an absence of data: a systematic review. *Global Health Action, 11*, 1-

901  14.

902  Celerier, C., Rouillon, I., Blanchard, M. et al. (2017). Pain and cochlear implantation: An unusual

903  complication? *Otol Neurotol, 38, 956-961.*

904  Cohen J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, N.J.:

905  Lawrence Erlbaum Associates.

906  Cohen, N.L., Waltzman, S.B., Fisher, S.G., et al. (1993). A prospective randomized study of cochlear

907  implants. *New England Journal of Medicine, 328,* 233-237.

908  Ching, T.Y.C., Incerti, P., Hill, M. (2004). Binaural benefits for adults who use hearing aids and

909  cochlear implants in opposite ears. *Ear Hear, 25, 9-21.*

910 Dakin, H., Devlin, N., Feng, Y., et al. (2015). The influence of cost-effectiveness and other factors on

911      NICE decisions. *Health Econ, 24,* 1256-1271.

912 Dawes, P., Munro, K.J. (2013). Unilateral and bilateral hearing aids, spatial release from masking and

913      auditory acclimatization. *J Acoust Soc Am., 134, 586-606.*

914 de Vet, H.C.W., Terwee, C.B., Knol, D.L., et al. (2006). When to use agreement versus reliability

915      measures. *J Clin Epid., 59, 1033-1039.*

916 Dolan, P. (1997). Modeling valuations for EuroQol health states. *Md Care, 35,* 1095-1108.

917 Dolan, P., Roberts, J. (2002). To what extent can we explain time trade-off values from other

918      information about respondents? *Social Science & Medicine, 54,* 919-929.

919 Drummond, M.F., O'Brien, B., Stoddart, G. L., et al. (2000). *Methods for the Economic Evaluation of*

920      *Health Care Programmes*, 2nd Edition. Oxford: Oxford University Press.

921 Dubois, R.W. (2016). Cost-effectiveness thresholds in the USA: are they coming? Are they already

922      here? *J Comp Eff Res, 5,* 9-11.

923 Feeny, D., Furlong, W., Torrance, G.W., et al. (2002). Multi-attribute and single-attribute utility

924      functions for the Health Utilities Index Mark 3 system. *Medical Care, 40,* 113-128.

925 Field, A. (2009). *Discovering Statistics Using SPSS,* 3rd Edition. London: Sage.

926 Finch, A.P., Brazier, J.E., Makuria, C., et al. (2017). An exploratory study on using principal-

927      component analysis and confirmatory factor analysis to identify bolt-on dimensions: the EQ-5D

928      case study. *Value in Health, 20,* 1362-1375.

929 Gatehouse, S., Noble, W. (2004). The Speech, Spatial and Qualities of Hearing Scale (SSQ). *Int J*

930      *Audiol, 43,* 85-99.

931   Gerlinger, C., Bamber, L., Leverkus, F., Schwenke, C., Haberland, C., Schmidt, G., Endikrat, J. (2019).

932   Comparing the EQ-5D-5L utility index based on value sets of different countries: impact on the

933   interpretation of clinical study results. *BMC Res Notes 12:18.*

934   Goman, A.M. (2014). *A comparison of bilateral cochlear implantation and bimodal aiding in severely-*

935   *profoundly hearing-impaired adults: head movements, clinical outcomes, and cost-effectiveness.*

936   Doctoral Thesis, University of York, UK.

937   Grutters, J.P.C., Joore, M.A., van der Horst, F., et al. (2007). Choosing between measures:

938   comparison of EQ-5D, HUI2 and HUI3 in persons with hearing complaints. *Quality of Life*

939   *Research, 16,* 1439-1449.

940   Härkönen, K., Kivekäs, I., Rautiainen, M. et al. (2015). Sequential bilateral cochlear implantation

941   improves working performance, quality of life, and quality of hearing. *Acta Oto-Laryngologica,*

942   *135, 440-446.*

943   Health Quality Ontario (2018). Bilateral cochlear implantation: a health technology assessment. *Ont*

944   *Health Technol Assess Ser [internet], 18,* 1-139.

945   Hinderink, J.B., Krabbe, P.F.M., van den Broek, P. (2000). Development and application of a health-

946   realted quality of life instrument for adults with cochlear implants: The Nijmegen Cochlear

947   Implant Questionnaire. *Otolaryngology – Head and Neck Surgery, 123, 756-765.*

948   Horsman, J., Furlong, W., Feeny, D., Torrance, G. (2003). The Health Utilities Index (HUI®): concepts,

949   measurement properties and applications. *Health Qual Life Outcomes, 1*, 1-13.

950   Hughes, K.C., Galvin, K.L. (2013). Measuring listening effort expended by adolescents and young

951   adults with unilateral or bilateral cochlear implants or normal hearing. *Cochlear Implants Int.,*

952   *14, 121-129.*

953   IBM Corp. (2019) *IBM SPSS Statistics for Windows, Version 26.0.* Armonk, NY: IBM Corp.

954    Jones, C. A., Feeny, D., Eng, K. (2005). Test-retest reliability of health utilities index scores: evidence

955        from hip fracture. *Int J Technol Assess Health Care, 21,* 393-398.

956    Kerber, S., Seeber, B.U. (2012). Sound localization in noise by normal-hearing listeners and cochlear

957        implant users. *Ear Hear, 33, 445-457.*

958    Kitterick, P.T., Lovett, R.E.S., Goman, A.M. et al. (2011). The AB-York crescent of sound: an apparatus

959        for assessing spatial-listening skills in children and adults. *Cochlear Implants Int., 12, 164-169.*

960    Koo, T.K., Li, M.Y. (2016). A guideline of selecting and reporting intraclass correlation coefficients for

961        reliability research. *Journal of Chiropractic Medicine*, *15, 155-163.*

962    Krabbe, P.F.M., Stouthard, M.E.A., Essink-Bot, M-L., et al. (1999). The effect of adding a cognitive

963        dimension to the EuroQol Multiattribute Health-status classification system. *J Clin Epidemiol,*

964        *52,* 293-301.

965    Kuthubutheen, J., Mittmann, N., Amoodi, H., et al. (2015). The effect of different utility measures on

966        the cost-effectiveness of bilateral cochlear implantation. *Laryngoscope, 125,* 442-447.

967    Lakens, D. (2017). Equivalence tests: a practical primer for t tests, correlations, and meta-analyses.

968        *Social Psychological and Personality Science* 8, 335-362.

969    Lammers, M.J.W., Grolman, W., Smulders, Y.E., et al. (2011). The cost-utility of bilateral cochlear

970        implantation: a systematic review. *Laryngoscope, 121,* 2604-2609.

971    Litovsky, R., Parkinson, A., Arcaroli, J. et al. (2006). Simultaneous bilateral cochlear implantation in

972        adults: a multicenter clinical study. *Ear Hear, 27, 714-731.*

973    Longworth, L., Yang, Y., Young, T., et al. (2014). Use of generic and condition-specific measures of

974        health-related quality of life in NICE decision-making: a systematic review, statistical modelling

975        and survey. *Health Technol Assess, 18,* 1-254.

976 Luo, N., Johnson, J.A., Shaw, J.W., et al. (2009). Relative efficiency of the EQ-5D, HUI2, and HUI3

977   index scores in measuring health burden of chronic medical conditions in a population health

978   survey in the United States. *Medial Care, 47,* 53-60.

979 Macran, S. (2003). Test-retest performance of EQ-5D. In R. Brooks, R. Rabin, F. de Charro (Eds.), *The*

980   *Measurement and Valuation of Health Status Using EQ-5D: A European Perspective* (pp. 43-54).

981   Dordrecht, NL: Kluwer Academic Publishers.

982 McGraw, K.O., Wong, S.P. (1996). Forming inferences about some intraclass correlation coefficients.

983   *Psychological Methods* 1, 30-46.

984 National Institute for Health and Care Excellence (2013). *Guide to the Methods of Technology*

985   *Appraisal 2013. Process and Methods Guides No. 9.* London: National Institute for Health and

986   Care Excellence.

987 Neumann, P.J., Cohen, J.T., Weinstein, M.C. (2014). Updating cost-effectiveness – the curious

988   resilience of the $50,000-per-QALY threshold. *New England Journal of Medicine, 371,* 796-797.

989 Noble, W., Gatehouse, S. (2006). Effects of bilateral versus unilateral hearing aid fitting on abilities

990   measured by the Speech, Spatial, and Qualities of Hearing Scale (SSQ). *Int J Audiol., 45, 172-*

991   *181.*

992 Noble, W., Tyler, R., Dunn, C., et al. (2008). Unilateral and bilateral cochlear implants and the

993   implant-plus-hearing-aid profile: comparing self-assessed and measured abilities. *Int J Audiol.,*

994   *47, 505-14.*

995 Office for National Statistics (2015). Period life expectancy by sex and 5-year age group, England,

996   Wales and regions of England 2011 to 13 and 2012 to 2014. Downloaded from

997   https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/lifeexpec

998   tancies/adhocs/005097periodlifeexpectancybysexand5yearagegroupenglandwalesandregionsof

999   england2011to13and2012to2014 on 10[th] September 2019.

1000     Oppe, M., Rand-Hendriksen, K., Shah, K., et al. (2016). EuroQol protocols for time trade-off valuation

1001        of health outcomes. *PharmacoEconomics, 34,* 993-1004.

1002     Palmer, C.S., Niparko, J.K., Wyatt, J.R., et al. (1999). A prospective study of the cost-utility of the

1003        multichannel cochlear implant. *Arch Otolaryngol Head Neck Surg, 125,* 1221-1229.

1004     Palta, M., Chen, H-Y., Kaplan, R.M., et al. (2011) Standard error of measurement of 5 health utility

1005        indexes across the range of health for use in estimating reliability and responsiveness. *Med*

1006        *Decis Making, 31, 260-269.*

1007     Potts, L.G., Skinner, M.W., Litovsky, R.A., et al. (2009). Recognition and localization of speech by

1008        adult cochlear implant recipients wearing a digital hearing aid in the nonimplanted ear (bimodal

1009        hearing). *J Am Acad Audiol., 20, 353-373.*

1010     Ramakers, G.G.J., van Zon, A., Stegeman, I., et al. (2015). The effect of cochlear implantation on

1011        tinnitus in patients with bilateral hearing loss: A systematic review. *Laryngoscope, 125, 2584-*

1012        *2592.*

1013     Rawlins, M.D., Culyer, A.J. (2004). National Institute for Clinical Excellence and its value judgements.

1014        *BMJ, 329,* 224-227.

1015     Rennies, J., Kidd, G. (2018). Benefit of binaural listening as revealed by speech intelligibility and

1016        listening effort. *J Acoust Soc Am., 144, 2147-2159.*

1017     Smulders, Y.E. (2016) *Unilateral versus simultaneous bilateral cochlear implantation in adults; a*

1018        *randomized controlled trial.* Doctoral Thesis, University of Utrecht, The Netherlands.

1019     Smulders, Y.E., van Zon, A., Stegeman, I., et al. (2016). Cost-utility of bilateral versus unilateral

1020        cochlear implantation in adults: a randomized controlled trial. *Otol Neurotol, 37,* 38-45.

1021  Summerfield, A.Q., Barton, G.R. (2019). Sensitivity of EQ-5D-3L, HUI2, HUI3, and SF-6D to changes in

1022      speech perception and tinnitus associated with cochlear implantation. *Qual Life Res, 28,* 1145-

1023      1154.

1024  Summerfield, A.Q., Barton, G.R., Toner, J., et al. (2006). Self-reported benefits from successive

1025      bilateral cochlear implantation in post-lingually deafened adults: randomised controlled trial.

1026      *International Journal of Audiology, 45(Supplement 1),* S99-S107.

1027  Summerfield, A.Q., Marshall, D.H., Barton, G.R., et al. (2002). A cost-utility scenario analysis of

1028      bilateral cochlear implantation. *Archives of Otolaryngology Head Neck Surgery*, *128,* 1255-1262.

1029  Swinburn, P., Lloyd, A., Boye, K.S., et al. (2013). Development of a disease-specific version of the EQ-

1030      5D-5L for use in patients suffering from psoriasis: lessons learned from a feasibility study in the

1031      UK. *Value in Health, 16,* 1156-1162.

1032  Theriou, C., Fielden, C.A., Kitterick, P.T. (2019). The cost-effectiveness of bimodal stimulation

1033      compared to unilateral and bilateral cochlear implant use in adults with bilateral severe to

1034      profound deafness. *Ear Hear, 40,* 1425-1436.

1035  UK Cochlear Implant Study Group (2004). Criteria of candidacy for unilateral cochlear implantation in

1036      postlingually deafened adults II: cost-effectiveness analysis. *Ear Hear, 25,* 336-360.

1037  Versteegh, M. M., Leunis, A., Uyl-de Groot, C. A., et al. (2012). Condition-specific preference-based

1038      measures: benefit or burden? *Value in Health, 15,* 504-513.

1039  Vickers, A.J., Altman, D.G. (2001). Analysing controlled trials with baseline and follow up

1040      measurements. *BMJ, 323,* 1123-1124.

1041  Vickers, D., De Raeve, L., Graham, J. (2016). International survey of cochlear implant candidacy.

1042      *Cochlear Implants International, 17:sup 1,* 36-41.

1043   World Health Organization (2001). *International Classification of Functioning Disability and Health.*

1044      WHO: Geneva.

1045   Yang, Y., Brazier, J.E., Tsuchiya, A., et al. (2011). Estimating a preference-based index for a 5-

1046      dimensional health state classification for asthma derived from the Asthma Quality of Life

1047      Questionnaire. *Med Decis Making, 31,* 281-291.

1048   Yang, Y., Rowan, D., Brazier, J., et al. (2015). An exploratory study to test the impact of three "bolt-

1049      on" items to the EQ-5D. *Value in Health, 18,* 522-60.

1050   York Health Economics Consortium (2016). Time Trade-Off. Accessed at

1051      https://yhec.co.uk/glossary/time-trade-off/ on 1st January 2021.

1052

1053   **Reference Note**

1054   Reference Note 1. On p. 513 of Versteegh et al (2012), the wording of the third criterion should be

1055   'shown to be superior to the existing generic PBM' rather than 'shown to be superior to the existing

1056   CS-PBM' (Matthias Versteegh, personal communication, 7th September 2018).

1057

1058   **Figure Legends**

1059   **Figure 1**  Relationship between mean utility and level of dimensions in Experiment 1a (upper panels)

1060   and 1b (lower panels). Each symbol plots the mean utility for one level of a dimension (filled

1061   symbols, non-students; open symbols, students).

1062   **Figure 2**  Influence of dimensions in Experiment 1a (upper panel) and 1b (lower panel). Filled bars,

1063   non-students; open bars, students. Error bars plot 95% confidence intervals of mean values.

1064   **Figure 3**  YBRQL Questionnaire. The preamble in this example is worded for users of bilateral

1065   cochlear implants.

1066 **Figure 4** Relationship between SSQ score and utility when participants considered their functional

1067 abilities using their normal configuration of devices (upper row of panels). Relationship between

1068 change in SSQ score and change in utility between monaural and binaural listening for members of

1069 the bimodal and bilateral groups (lower row of panels). (Panels A and E, YBRQL; B and F, HUI3

1070 Hearing*;* D and G, HUI3 Health*;* D and H, EQ-5D-3L). (Black circles, Unilateral Group; dark grey

1071 squares, Bimodal Group; light grey triangles, Bilateral Group).

1072 **Figure 5** Box plots of SSQ score (Panels A and C) and utility (Panels B and D) when the Bilateral

1073 Group (Panels A and B) and the Bimodal Group (Panels C and D) considered their functioning with

1074 their first or only implant (open boxes) and with two devices (filled boxes). Whiskers mark 10th and

1075 90th percentiles; box marks 25th and 75th percentiles; heavy line within box marks 50th percentile;

1076 filled circles plot outliers beyond the 10th and 90th percentiles. (Panel E) z scores, significance levels

1077 (p), and effect sizes (R) from Wilcoxon signed ranks tests comparing scores with 1 and 2 devices.

1078 **Figure 6** Box plots of levels of function on each dimension of the YBRQL when the Bimodal and

1079 Bilateral Groups considered their function using their first or only implant (open boxes) or two

1080 devices (filled boxes). Results of Wilcoxon signed ranks tests comparing reported levels with one and

1081 two devices (*, p<.05; **, p<.01).

1082 **Figure 7** Decision matrix for interpreting binaural incremental cost-effectiveness ratios (ICERs)

1083 informed by the YBRQL together with health ICERs informed by a generic PBM.

1084

1085 **Table Legends**

1086 **Table 1** Estimates of the gain in health utility between unilateral and bilateral cochlear implantation

1087 in bilaterally-impaired post-lingually deafened adults obtained with the EuroQol Descriptive System

1088 (EQ-5D-3L), the Health Utilities Index Mark III (HUI3), and versions of the Time Trade-off technique

1089 (TTO).

1090     **Table 2**  Vignettes describing levels of difficulty on three dimensions of binaural hearing.

1091     **Table 3**  Numbers (N) and age of participants in Experiments 1a and 1b.

1092     **Table 4**  Calculation of age-standardized binaural utilities.

1093     **Table 5** YBRQL value set. Entries are the 125 values of the expanded value set. Underlined entries

1094     are the 27 values of the initial value set.

1095     **Table 6**  Kendall's coefficients of correlation among outcome measures (N=28). (*$p$<.05, **$p$<.01).

1096     **Table 7**  Discriminative ability: Measures of binaural advantage for the Bimodal and Bilateral Groups,

1097     and of the difference in binaural advantage (Bilateral – Bimodal), with 95% confidence intervals (95%

1098     CI) estimated by bootstrapping (3,000 samples per estimate, bias-corrected and accelerated).

1099     **Table 8** Estimates of the loss of utility due to impaired hearing and the gain in utility from using a

1100     second device estimated. (IQR: inter-quartile range.)

1101

1102    **Supplementary Digital Content**

1103    1. SDC 1 (Gain in HRQL associated with bilateral cochlear implantation).pdf

1104    2. SDC 2 (Levels of function reported with the HUI3 and the YBRQL).pdf

1105    3. SDC 3 (Response booklet for Experiment 1a).pdf

1106    4. SDC 4 (Response booklet for Experiment 1b).pdf

1107    5. SDC 5 (Pilot Experiment and test of equivalence).pdf

1108    6. SDC 6 (Response booklet for Pilot Experiment).pdf

1109    7. SDC 7 (Relationship between age and binaural utility).pdf

1110    8. SDC 8 (Tests of reproducibility).pdf

1111    9. SDC 9 (Editable version of the YBRQL questionnaire).docx

1112

Table 1  Estimates of the gain in health utility between unilateral and bilateral cochlear implantation in bilaterally-impaired post-lingually deafened adults obtained with the EuroQol Descriptive System (EQ-5D-3L), the Health Utilities Index Mark 3 (HUI3), and versions of the Time Trade-off technique (TTO).

| Measure | Study | Design | Number of participants | Mean increment in utility (95% CI)[a] | Significance of increment |
|---|---|---|---|---|---|
| EQ-5D-3L | Summerfield et al. (2006) | Randomised controlled trial[b] | 2x12 | -.006 (-.091 to .078)[c] | n.s. |
| | | | 24 | -.063(-.120 to .005)[d] | n.s. |
| | Kuthubutheen et al. (2014) | Scenario analysis[e] | 142 | .04 | p<.05 |
| | Smulders et al. (2016) | Randomised controlled trial[f] | 2x19 | -.02 | n.s. |
| HUI3 | Summerfield et al. (2006) | Randomised controlled trial[b] | 2x12 | .105 (-.073 to .282)[c] | n.s. |
| | | | 24 | -.015 (-.110 to .079)[d] | n.s. |
| | | | 24 | .030 (-.045 to .104)[g] | n.s. |
| | Bichey & Miyamoto (2008) | Retrospective study[h] | 23 | .12 (.09 to .14) | p<.001 |
| | Kuthubutheen et al. (2014) | Scenario analysis[e] | 142 | .035 | p<.05 |
| | Smulders et al. (2016) | Randomised controlled trial[f] | 2 x 19 | .04 | n.s. |
| TTO | Summerfield et al. (2002) | Scenario analysis[i] | 70 | .03 (.02 to .04) | p<.05 |
| | Kuthubutheen et al. (2014) | Scenario analysis[j] | 142 | .12 | p<.05 |
| | Smulders et al. (2016) | Randomised controlled trial[k] | 2x19 | .09 | n.s. |

[a]CI = Confidence Interval where reported or calculable from published data.

[b]24 users of one implant were randomised to receive a second implant either immediately (treatment group, N=12) or after 12 months (control group, N=12).

[c]Comparison of groups 9 months after treatment group received second implant.

[d]Before-and-after comparison of combined group 9 months after each patient had received second implant.

[e]Scenario analysis of unilateral and bilateral implantation; informants were candidates for implantation (N=30), users of one implant (N=30), users of two implants (N=30), and clinicians (N=52).

[f]38 candidates for implantation were randomised to receive either bilateral implantation (treatment group, N=19) or unilateral implantation (control group, N=19).

[g]Before-and-after comparison of combined group 9 months after each patient had received second implant with control for changes in annoyance due to tinnitus.

[h]Retrospective assessment of health status with one implant, and contemporary assessment of health status with two implants, by users of bilateral implants.

[i]Scenario analysis of unilateral and bilateral implantation with researchers and clinicians as informants. Time frame: current age to 75 years.

[j]Scenario analysis of unilateral and bilateral implantation; informants were candidates for implantation (N=30), users of one implant (N=30), users of two implants (N=30), and clinicians (N=52). Time frame: 30 years.

[k]38 candidates for implantation were randomised to receive either bilateral implantation (treatment group, N=19) or unilateral implantation (control group, N=19). Time frame: life expectancy.

Table 2  Vignettes describing levels of difficulty on three dimensions of binaural hearing.

| Dimension | Level | Vignette |
|---|---|---|
| Speech Perception in Noise | 1: No difficulty | When a friend speaks to you while the TV is on or other people are chatting in the same room, you can hear your friend speaking easily, usually picking up all of the words they say. |
| | 2: Some difficulty | When a friend speaks to you while the TV is on or other people are chatting in the same room, you can hear your friend speaking, but you can only pick out some of the words they say. This can lead to confusion if you miss an important word. Sometimes you need them to repeat themselves or to turn the volume down for you to understand them. |
| | 3: Great difficulty | When a friend speaks to you while the TV is on or other people are chatting in the same room, you find it very difficult to hear your friend speaking. You are usually unable to pick out the words they say. This regularly leads to misunderstanding and confusion. The room needs to be completely quiet for you to understand them. |
| Localisation | 1: No difficulty | You can work out where sounds are coming from accurately. You can point to where a sound is coming from easily. |
| | 2: Some difficulty | You have some difficulty working out where sounds are coming from. You can usually tell if a sound is coming from the right- or left-hand side, but you cannot be more accurate than that. As a result, you are not always sure who is speaking when you are in a group with several people. |
| | 3: Great difficulty | You have great difficulty working out where sounds are coming from. You cannot even tell if a sound is coming from the right- or left-hand side without looking around. As a result, you find it very difficult to tell who is speaking when you are in a group with several people. You are also worried about your safety outdoors because of your difficulty working out where sounds are coming from. |
| Effort and Fatigue | 1: No difficulty | You have to concentrate a little when you are trying to hear something or someone. You can hear what people are saying with only a little effort.  By the end of the day, you are not mentally or physically tired because of your hearing. |
| | 2: Some difficulty | You have to concentrate quite hard when you are trying to hear something or someone. You have to put in some effort to hear what people are saying. By the end of the day, you are moderately mentally and physically tired because of your hearing. |
| | 3: Great difficulty | You have to concentrate very hard when you are trying to hear something or someone. You have to put in a great deal of effort to hear what people are saying. By the end of the day, you are extremely mentally and physically tired because of your hearing. |

Table 3   Numbers (N) and age of participants in Experiments 1a and 1b

| Experiment | Group | Participants (N) | Inconsistent traders (N) | Zero traders (N) | Included in analyses (N) | Minimum age (years) | Mean age (years) | Maximum age (years) | % female |
|---|---|---|---|---|---|---|---|---|---|
| 1a | Students | 59 | 0 | 0 | 59 | 18 | 20.8 | 26 | 79.7 |
|  | Non-students [Public] | 52 | 1 | 4 | 51 | 22 | 44.1 | 64 | 52.9 |
| 1b | Students | 48 | 0 | 0 | 48 | 18 | 20.7 | 26 | 60.9 |
|  | Non-students [Public] | 48 | 3 | 5 | 45 | 18 | 49.1 | 80 | 57.8 |

Table 4  Calculation of age-standardized binaural utilities.

| Age Group (years) | Number in adult population | Proportion in adult population | Students N | Students Mean | Non-Students N | Non-Students Mean | Combined group N | Combined group Mean | Age-weighted contribution |
|---|---|---|---|---|---|---|---|---|---|
| | UK population data [a] | | | | Binaural utility data [b] | | | | |
| 18-20 | 2,390,852 | 0.047 | 52 | 0.748 | | | 52 | 0.748 | 0.035 |
| 21-30 | 8,815,943 | 0.172 | 55 | 0.736 | 15 | 0.730 | 70 | 0.735 | 0.126 |
| 31-40 | 8,379,311 | 0.163 | | | 14 | 0.827 | 14 | 0.827 | 0.135 |
| 41-50 | 9,070,302 | 0.177 | | | 27 | 0.847 | 27 | 0.847 | 0.150 |
| 51-60 | 8,282,160 | 0.161 | | | 31 | 0.868 | 31 | 0.868 | 0.140 |
| >60 | 14,400,593 | 0.281 | | | 9 | 0.792 | 9 | 0.792 | 0.222 |
| Total | 51,339,161 | 1.000 | 107 | | 96 | | 203 | | 0.808 |

[a] Office for National Statistics (2015); [b] From Experiments 1a and 1b.

Table 5  YBRQL value set. Entries are the 125 values of the expanded value set. Underlined entries are the 27 values of the initial value set.

| Speech-in-noise Level | Localisation Level | Effort & Fatigue Level | | | | |
|---|---|---|---|---|---|---|
| | | 1 | 1.5 | 2 | 2.5 | 3 |
| 1 | 1 | 0.9625 | 0.9272 | 0.8919 | 0.8556 | 0.8192 |
| 1 | 1.5 | 0.9323 | 0.8978 | 0.8632 | 0.8359 | 0.8087 |
| 1 | 2 | 0.9021 | 0.8683 | 0.8345 | 0.8163 | 0.7981 |
| 1 | 2.5 | 0.8758 | 0.8491 | 0.8223 | 0.7949 | 0.7674 |
| 1 | 3 | 0.8495 | 0.8298 | 0.8101 | 0.7734 | 0.7367 |
| | | | | | | |
| 1.5 | 1 | 0.9321 | 0.9001 | 0.8681 | 0.8320 | 0.7958 |
| 1.5 | 1.5 | 0.9034 | 0.8743 | 0.8451 | 0.8147 | 0.7843 |
| 1.5 | 2 | 0.8746 | 0.8484 | 0.8222 | 0.7974 | 0.7727 |
| 1.5 | 2.5 | 0.8526 | 0.8291 | 0.8055 | 0.7768 | 0.7481 |
| 1.5 | 3 | 0.8305 | 0.8097 | 0.7888 | 0.7562 | 0.7235 |
| | | | | | | |
| 2 | 1 | 0.9017 | 0.8730 | 0.8443 | 0.8084 | 0.7725 |
| 2 | 1.5 | 0.8744 | 0.8508 | 0.8271 | 0.7935 | 0.7599 |
| 2 | 2 | 0.8472 | 0.8285 | 0.8098 | 0.7786 | 0.7473 |
| 2 | 2.5 | 0.8294 | 0.8090 | 0.7887 | 0.7588 | 0.7288 |
| 2 | 3 | 0.8116 | 0.7896 | 0.7676 | 0.7389 | 0.7103 |
| | | | | | | |
| 2.5 | 1 | 0.8889 | 0.8572 | 0.8256 | 0.7950 | 0.7644 |
| 2.5 | 1.5 | 0.8646 | 0.8390 | 0.8135 | 0.7803 | 0.7472 |
| 2.5 | 2 | 0.8403 | 0.8208 | 0.8014 | 0.7657 | 0.7301 |
| 2.5 | 2.5 | 0.8234 | 0.7989 | 0.7745 | 0.7442 | 0.7140 |
| 2.5 | 3 | 0.8065 | 0.7770 | 0.7476 | 0.7227 | 0.6979 |
| | | | | | | |
| 3 | 1 | 0.8761 | 0.8415 | 0.8069 | 0.7816 | 0.7563 |
| 3 | 1.5 | 0.8548 | 0.8273 | 0.7999 | 0.7672 | 0.7346 |
| 3 | 2 | 0.8335 | 0.8132 | 0.7929 | 0.7529 | 0.7128 |
| 3 | 2.5 | 0.8174 | 0.7888 | 0.7602 | 0.7297 | 0.6991 |
| 3 | 3 | 0.8014 | 0.7645 | 0.7276 | 0.7065 | 0.6854 |

Table 6  Kendall's coefficients of correlation among outcome measures (N=28). (*$p$<.05, **$p$<.01).

| | YBRQL (binaural utility) | HUI3 (hearing utility) | HUI3 (health utility) | EQ-5D-3L (health utility) |
|---|---|---|---|---|
| Localisation | .49** | .20 | .27 | .17 |
| Speech in noise | -.30* | -.25 | -.28* | -.09 |

Table 7  Discriminative ability: Measures of binaural advantage for the Bimodal and Bilateral Groups, and of the difference in binaural advantage (Bilateral – Bimodal), with 95% confidence intervals (95% CI) estimated by bootstrapping (3,000 samples per estimate, bias-corrected and accelerated).

| | Group | | |
|---|---|---|---|
| | Bimodal | Bilateral | Bilateral - Bimodal |
| Measure | Binaural advantage (95% CI) | Binaural advantage (95% CI) | Difference in binaural advantage (95% CI)[a] |
| SSQ score | .712 (.207 to 1.142)[b] | 2.866 (1.679 to 4.156)[b] | 2.154 (.749 to 3.842)[b] |
| YBRQL (Binaural utility) | .031 (.009 to .070)[b] | .098 (.067 to .127)[b] | .067 (.010 to .111)[b] |
| HUI3 (Hearing utility) | .102 (.018 to .183)[b] | .123 (.041 to .213)[b] | .021 (-.074 to .117) |
| HUI3 (Health utility) | .037 (-.017 to .087) | .065 (.022 to .122)[b] | .029 (-.050 to .112) |
| EQ-5D-3L (Health utility) | .024 (-.030 to .088) | .009 (-.012 to .037) | -.015 (-.100 to .056) |

[a] Differences estimated by bootstrapping need not exactly equal the difference between the estimates of the Bimodal and Bilateral means.

[b] 95% confidence interval does not include zero.

Table 8  Estimates of the loss of utility due to impaired hearing and the gain in utility from using a second device. (IQR: inter-quartile range.)

| Participants | N | Measure | Median | IQR |
|---|---|---|---|---|
| | | Loss of utility due to impaired hearing when using first or only implant | | |
| All Groups | 28 | YBRQL | .174 | .131 to .216 |
| | | HUI3 | .148 | .142 to .257 |

| Participants | N | Measure | Median | IQR |
|---|---|---|---|---|
| | | Gain in utility from using a second device | | |
| Bimodal Group | 9 | YBRQL | .020 | .000 to .046 |
| | | HUI3 | .000 | -.046 to .191 |
| | | EQ-5D-3L | .000 | .000 to .075 |
| Bilateral Group | 11 | YBRQL | .107 | .035 to .162 |
| | | HUI3 | .000 | .000 to .170 |
| | | EQ-5D-3L | .000 | .000 to .000 |

Summerfield Kitterick Goman (Revised)

This document contains the figure legends followed by the figures of the revised paper.

Figure 1  Relationship between mean utility and level of dimensions in Experiment 1a (upper panels) and 1b (lower panels). Each symbol plots the mean utility for one level of a dimension (filled symbols, non-students; open symbols, students).

Figure 2  Influence of dimensions in Experiment 1a (upper panel) and 1b (lower panel). Filled bars, non-students; open bars, students. Error bars plot 95% confidence intervals of mean values.

Figure 3  YBRQL Questionnaire. The preamble in this example is worded for users of bilateral cochlear implants.

Figure 4  Relationship between SSQ score and utility when participants considered their functional abilities using their normal configuration of devices (upper row of panels).  Relationship between change in SSQ score and change in utility between monaural and binaural listening for members of the bimodal and bilateral groups (lower row of panels). (Panels A and E, YBRQL; B and F, HUI3 Hearing; D and G, HUI3 Health; D and H, EQ-5D-3L). (Black circles, Unilateral Group; dark grey squares, Bimodal Group; light grey triangles, Bilateral Group).

Figure 5  Box plots of SSQ score (Panels A and C) and utility (Panels B and D) when the Bilateral Group (Panels A and B) and the Bimodal Group (Panels C and D) considered their functioning with their first or only implant (open boxes) and with two devices (filled boxes). Whiskers mark 10$^{th}$ and 90$^{th}$ percentiles; box marks 25$^{th}$ and 75$^{th}$ percentiles; heavy line within box marks 50$^{th}$ percentile; filled circles plot outliers beyond the 10$^{th}$ and 90$^{th}$ percentiles. (Panel E) z scores, significance levels (p), and effect sizes (R) from Wilcoxon signed ranks tests comparing scores with 1 and 2 devices.

Figure 6  Box plots of reported levels of function on each dimension of the YBRQL when thebilateral and bimodal groups considered their function using their first or only implant (open boxes) or two devices (filled boxes). Whiskers mark 10$^{th}$ and 90$^{th}$ percentiles; box marks 25$^{th}$ and 75$^{th}$ percentiles; heavy line within box marks 50$^{th}$ percentile; filled circle plots outlier beyond the 10$^{th}$ percentile. Results of Wilcoxon signed ranks tests comparing reported levels with one and two devices (✶, p<.05; ✶✶, p<.01).

Figure 7  Decision matrix for interpreting Binaural ICERs informed by the YBRQL together with Health ICERs informed by a generic PBM.
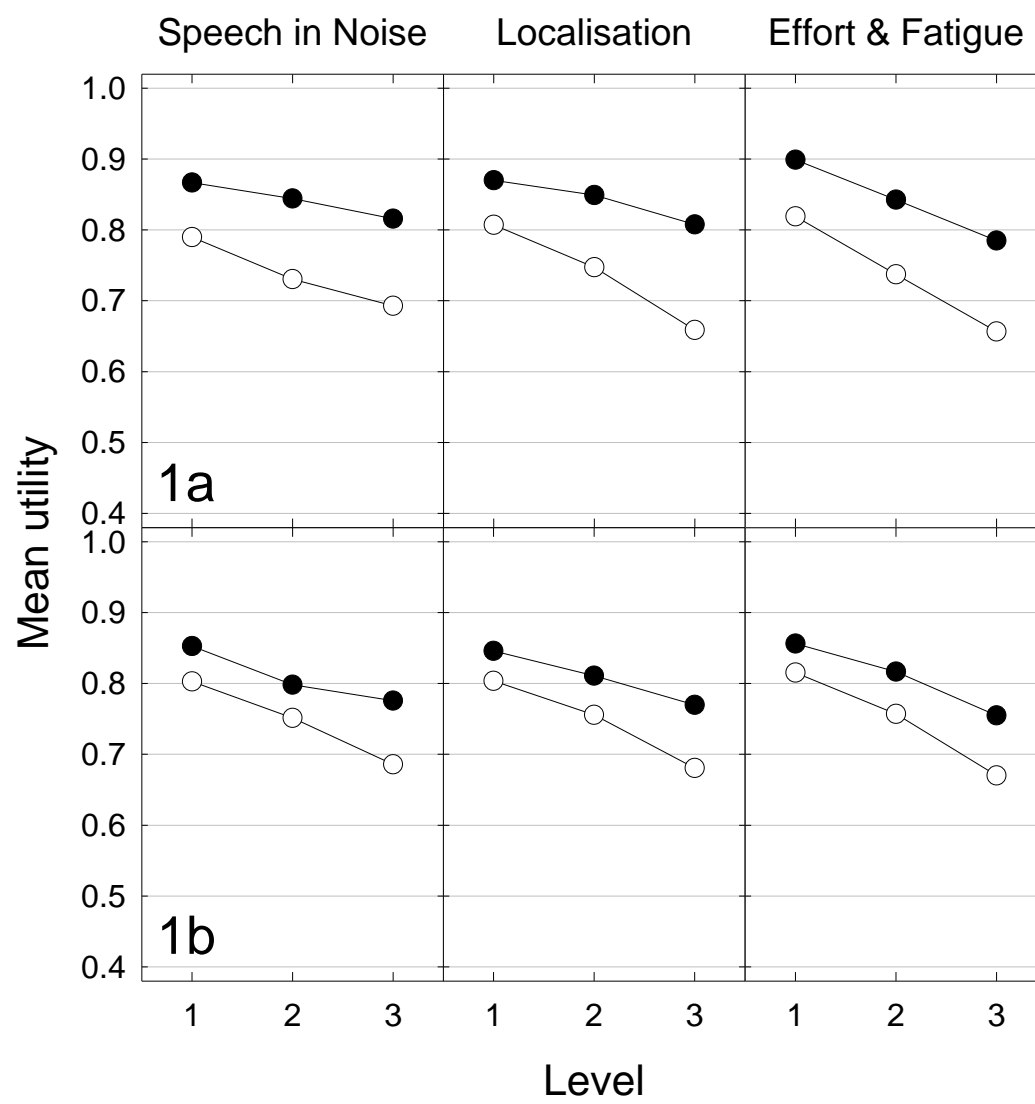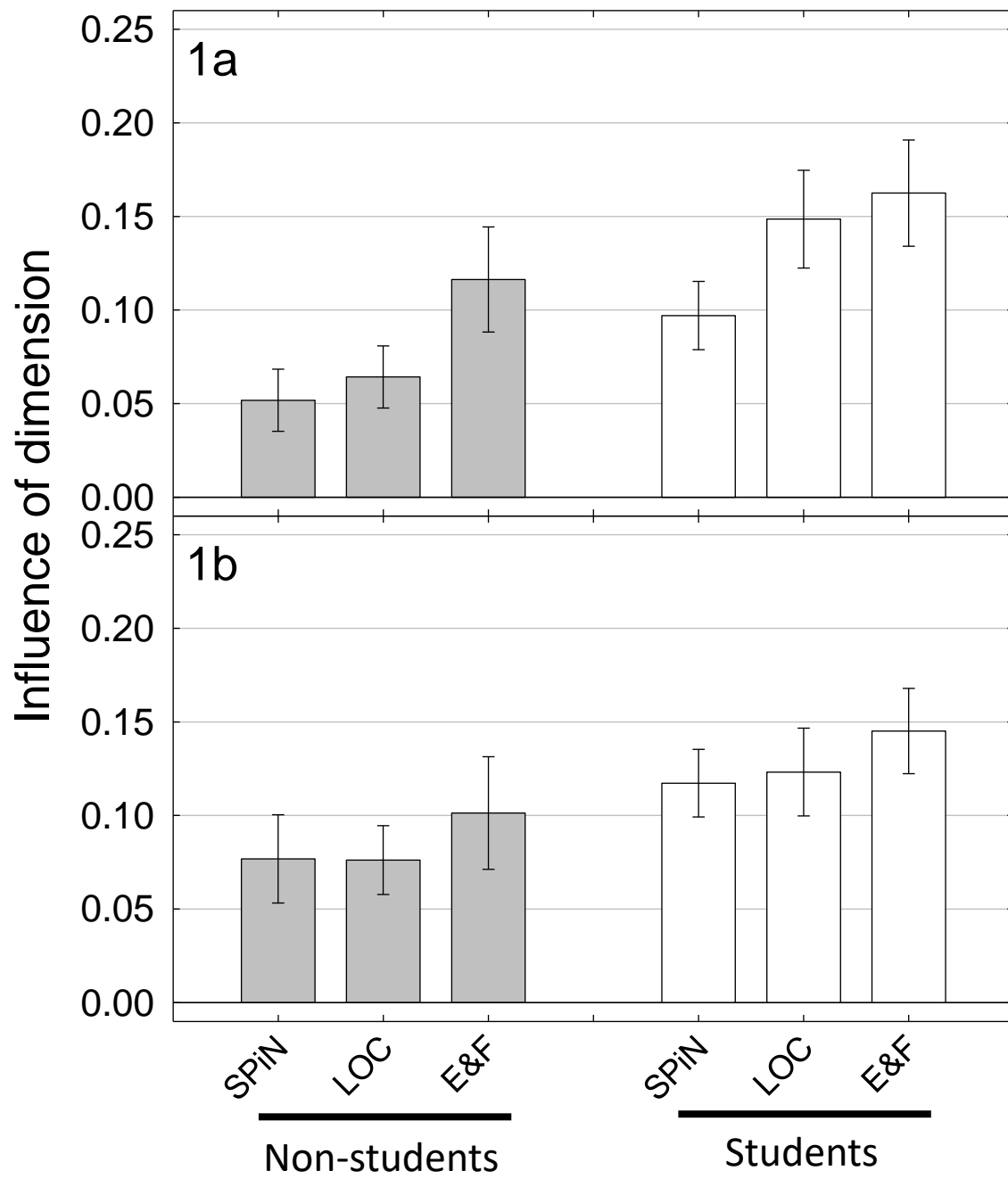
Figure 1

Figure 2

Figure 3

By placing a tick in one box under each heading below, please indicate which statement best describes your own hearing **when using your two cochlear implants together.**

Please do not tick more than one box under each heading.

Date:

P/N:

## Understanding speech when there is background noise

☐ 1. When a friend speaks to you while the TV is on or other people are chatting in the same room, you can hear your friend speaking easily, usually picking up all of the words they say.

☐ 2. Between 1 and 3.

☐ 3. When a friend speaks to you while the TV is on or other people are chatting in the same room, you can hear your friend speaking, but you can only pick out some of the words they say. This can lead to confusion if you miss an important word. Sometimes you need them to repeat themselves or to turn the volume down for you to understand them.

☐ 4. Between 3 and 5.

☐ 5. When a friend speaks to you while the TV is on or other people are chatting in the same room, you find it very difficult to hear your friend speaking. You are usually unable to pick out the words they say. This regularly leads to misunderstanding and confusion. The room needs to be completely quiet for you to understand them.

## Working out where sounds are coming from

☐ 1. You can work out where sounds are coming from accurately. You can point to where a sound is coming from easily.

☐ 2. Between 1 and 3.

☐ 3. You have some difficulty working out where sounds are coming from. You can usually tell if a sound is coming from the right- or left-hand side, but you cannot be more accurate than that. As a result, you are not always sure who is speaking when you are in a group with several people.

☐ 4. Between 3 and 5.

☐ 5. You have great difficulty working out where sounds are coming from. You cannot even tell if a sound is coming from the right- or left-hand side without looking around. As a result, you find it very difficult to tell who is speaking when you are in a group with several people. You are also worried about your safety outdoors because of your difficulty working out where sounds are coming from.

## Effort and fatigue

☐ 1. You have to concentrate a little when you are trying to hear something or someone. You can hear what people are saying with only a little effort. By the end of the day, you are **not** mentally or physically tired because of your hearing.

☐ 2. Between 1 and 3.

☐ 3. You have to concentrate quite hard when you are trying to hear something or someone. You have to put in some effort to hear what people are saying. By the end of the day, you are moderately mentally and physically tired because of your hearing.

☐ 4. Between 3 and 5.

☐ 5. You have to concentrate very hard when you are trying to hear something or someone. You have to put in a great deal of effort to hear what people are saying. By the end of the day, you are extremely mentally and physically tired because of your hearing.
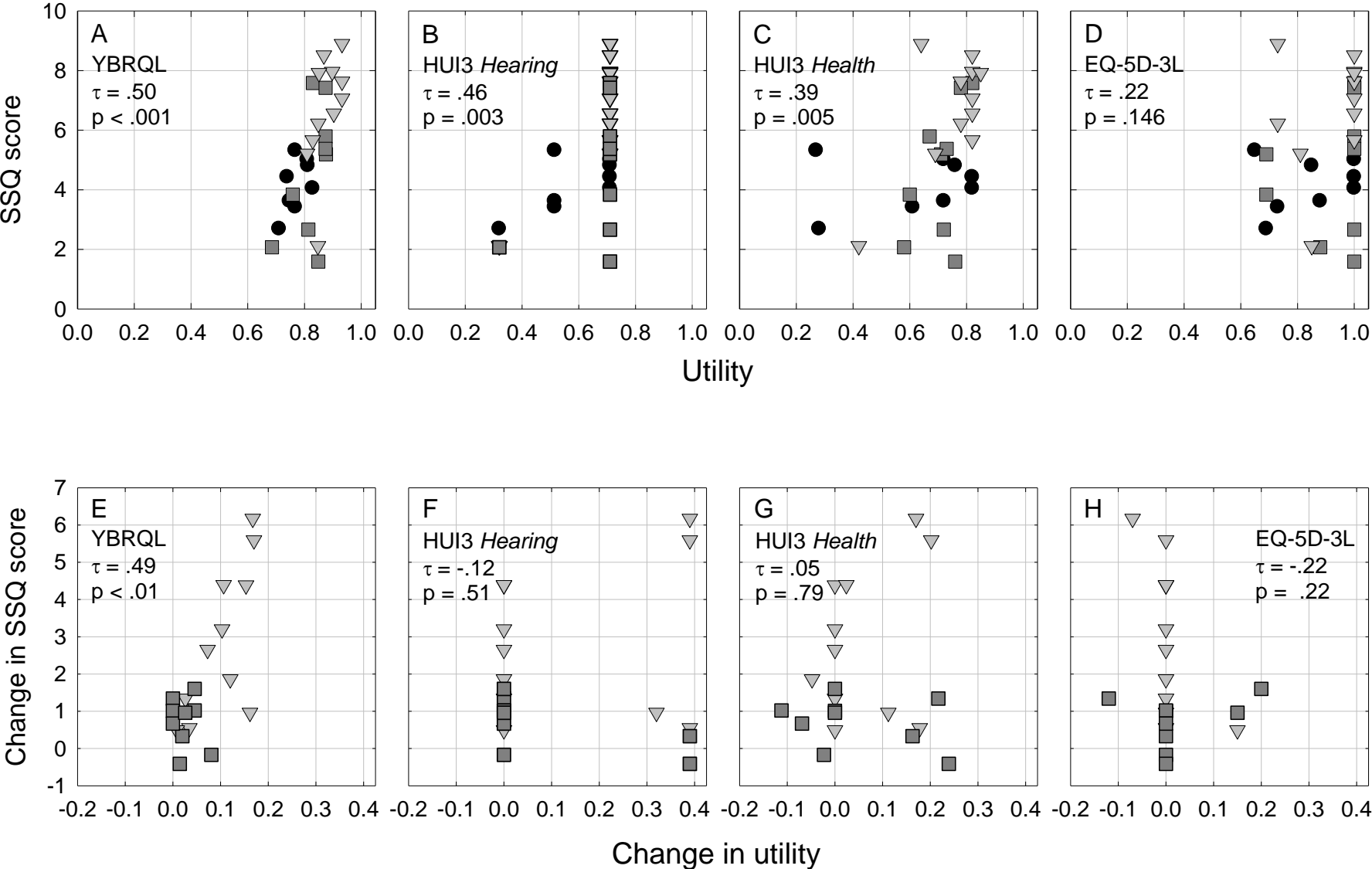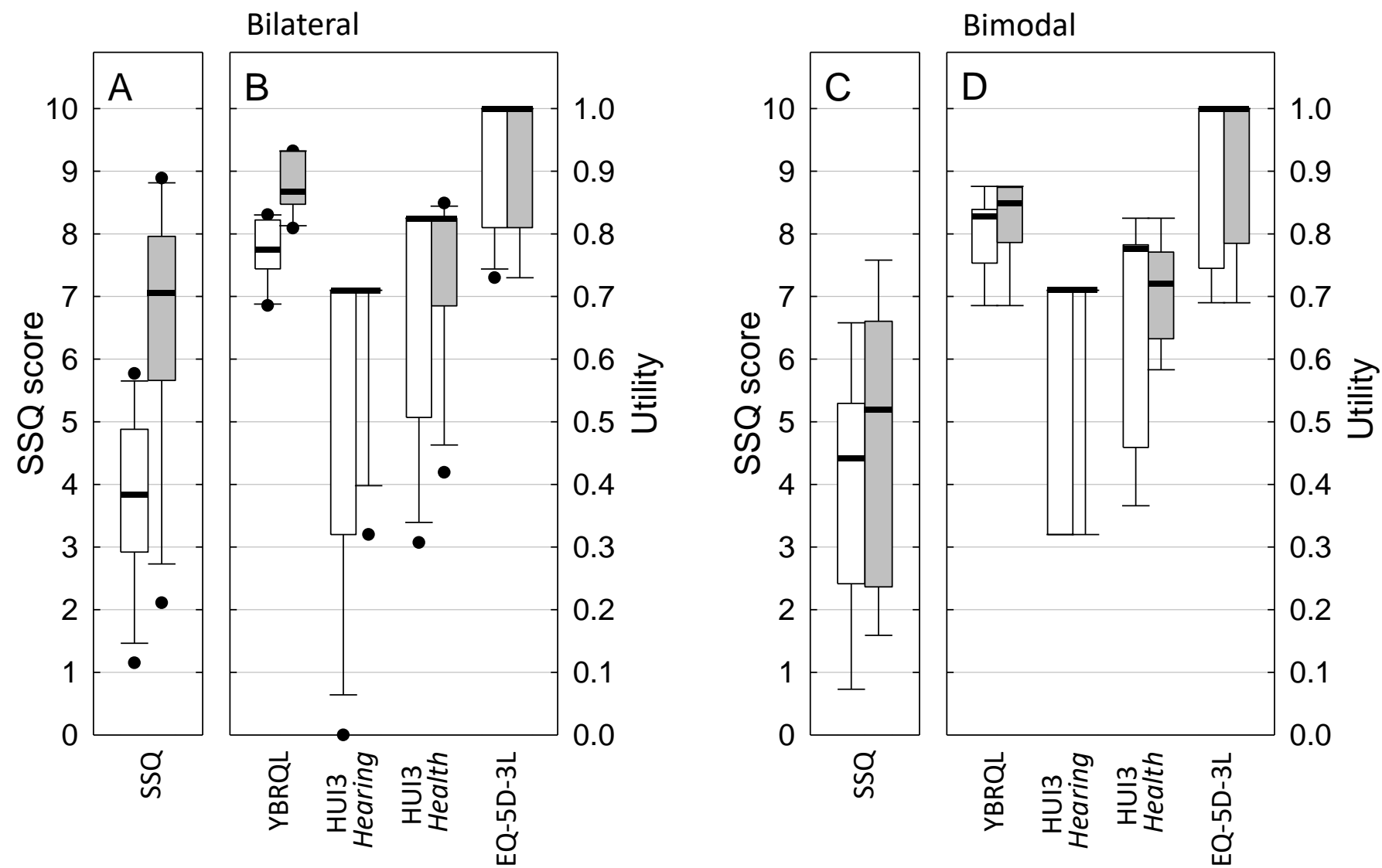
Figure 4

Figure 5



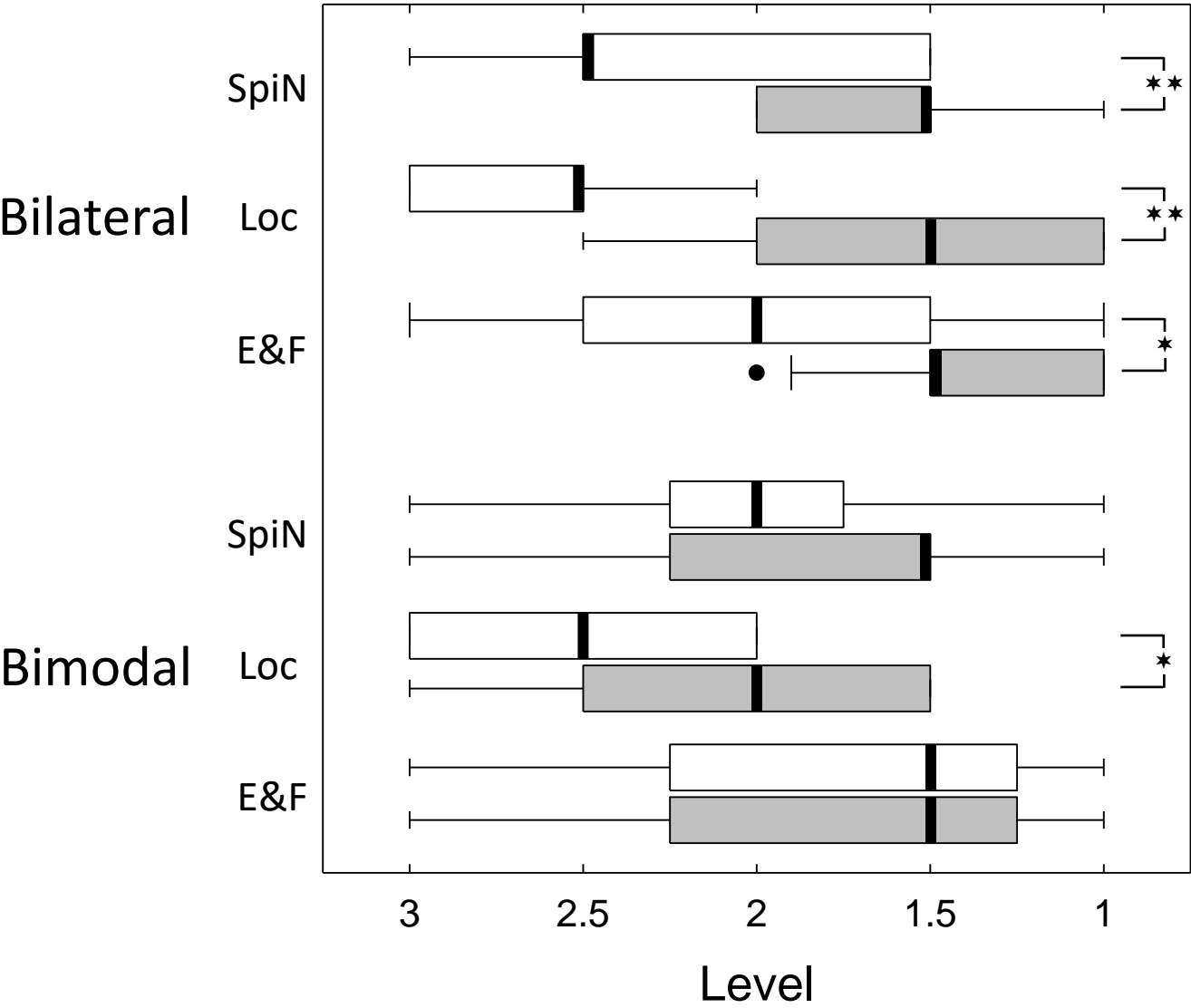| **E** | z | 2.93 | | 2.93 | 1.89 | 1.78 | .447 | | 2.19 | | 2.21 | 1.41 | .943 | 1.07 |
| | p | .003 | | .003 | .059 | .075 | .655 | | .028 | | .027 | .157 | .345 | .285 |
| | R | .63 | | .63 | .40 | .38 | .10 | | .52 | | .52 | .33 | .22 | .25 |

Figure 6

Figure 7

|  |  | Binaural | | |
|  |  | Binaural ICER is more favourable than WTP threshold | Binaural ICER is in the range of uncertainty | Binaural ICER is less favourable than the upper boundary of uncertainty |
| **Health** | Health ICER is more favourable than WTP Threshold | 1. *Treatment is cost-effective.* | 2. *Treatment is cost-effective.* Investigate reasons why gain in binaural-related QoL is smaller than gain in health-related QoL. | 3. Investigate reasons for small or negative gain in binaural-related QoL before deciding whether treatment is cost-effective. |
| | Health ICER is in the range of uncertainty | 4. *Treatment is probably cost-effective.* | 5. *Treatment is unlikely to be cost-effective.* Investigate reasons why costs are high and/or gain in binaural-related QoL is small. | 6. *Treatment is unlikely to be cost-effective.* Investigate reasons why costs are high and/or gain in binaural-related QoL is small or negative. |
| | Health ICER is less favourable than the upper boundary of uncertainty | 7. *Treatment is not cost-effective.* Investigate reasons why costs are high and/or gain in health-related QoL is small or negative. | 8. *Treatment is not cost-effective.* Investigate reasons why costs are high and/or gain in binaural-related QoL is small. | 9. *Treatment is not cost-effective.* Investigate reasons why costs are high and/or gain in binaural-related QoL is small or negative. |