# Objective measures for detecting the Auditory Brainstem Response: comparisons of specificity, sensitivity, and detection time**.

Chesnaye M.A[1], Bell S.L[2], Harte J.M[3], & Simpson D.M[4].

1. Institute of Sound and Vibration Research, Faculty of Engineering and the Environment, University of Southampton, United Kingdom. mac1f14@soton.ac.uk

2. Institute of Sound and Vibration Research, Faculty of Engineering and the Environment, University of Southampton, United Kingdom. s.l.bell@soton.ac.uk

3. Interacoustics Research Unit, c/o Technical University of Denmark, Denmark. jmha@iru.interacoustics.com

4. Institute of Sound and Vibration Research, Faculty of Engineering and the Environment, University of Southampton, United Kingdom. ds@isvr.soton.ac.uk

** PLEASE NOTE : THIS DOCUMENT WAS THE FINAL SUBMISSION TO THE JOURNAL AND DOES NOT INCLUDE FINAL PROOF READING AND FORMATTING. FOR THOSE WITH ACCESS, THE FINAL COPYRIGHTED VERSION CAN BE OBTAINED ONLINE FROM THE INTERNATIONAL JOURNAL OF AUDIOLOGY

24    **Abstract**

25    **Objective:** To evaluate and compare the specificity, sensitivity, and detection time of various time-

26    domain and multi-band frequency domain methods when detecting the auditory brainstem response

27    (ABR).

28    **Design:** Simulations and subject recorded data were used to assess and compare the performance of

29    the Hotelling's $T^2$ test (applied in either time or frequency domain), two versions of the modified q-

30    sample uniform scores test, and both the Fsp and Fmp, which were evaluated using both conventional

31    F-distributions with assumed degrees of freedom and a bootstrap approach.

32    **Study Sample:** Data consisted of simulations along with click-evoked ABRs and recordings of EEG

33    background activity from 12 and 17 normal hearing adults respectively.

34    **Results**: An overall advantage in sensitivity and detection time was demonstrated for the Hotelling's

35    $T^2$ test. The false-positive rates (FPRs) of the Fsp and Fmp were also closer to the nominal alpha level

36    when evaluating statistical significance using the bootstrap approach, as opposed to using

37    conventional F-distributions. The FPRs of the remaining methods were slightly higher than expected.

38    **Conclusions:** In the current work, Hotelling's $T^2$ outperformed the alternative methods for

39    automatically detecting ABRs. Its promise as a sensitive and efficient detection method should now

40    be tested in a larger clinical study.

41

42

43

44

45

46

47

48 **Keywords**

49 Auditory Brainstem Response detection, Hotelling's $T^2$, Fsp, Fmp, Bootstrapping, modified q-sample
50 uniform scores test

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

# 1. Introduction

Transient auditory brainstem responses (ABRs) are defined as short changes in neural activity along the auditory pathway in response to a brief acoustic stimulus, such as a click, chirp, or tone burst. They are typically recorded non-invasively using surface mounted electrodes, and are used primarily for diagnosing abnormalities within the auditory system, such as hearing loss and various neurological disorders. Determining whether an ABR is present (by either inspecting the data visually, or by applying an objective statistical test) is usually the first step for these applications, after which additional analysis can be performed on, for example, the morphology of the response.

The focus for this paper is on objective methods for detecting the ABR. The goal is to compare the performance of various statistical detection methods in terms of (i) true-positive rates (TPRs – in the current paper defined as the fraction of ABR responses that is detected), (ii) false-positive rates (FPRs – defined as the fraction of cases with no response that were incorrectly deemed to have a response), and (iii) detection time, i.e. the number of stimuli (expressed in time) required for detecting a significant response, which can be considered as the three most important properties for ABR detection methods.

Most ABR detection methods are classified as either time or frequency domain techniques. In the frequency domain detection is more challenging due to the spectral content of the ABR being spread across multiple bands (Elberling, 1976; Kevanishvili & Aphonchenko, 1979; Elberling, 1979; Suzuki et al., 1982). As most frequency domain techniques are applied to a single spectral band, they would need to be applied multiple times to cover the bandwidth of a typical ABR. The latter can result in an inflated FPR, and adjusted critical decision boundaries are required in order to preserve the desired alpha level of the test. This process may result in a significantly lower test sensitivity.

The broadband spectral content of the ABR has therefore led the majority of scientific investigations to explore methods for assessing multiple spectral bands within a single test, i.e. multi-band detection

98    methods. A powerful multi-band detection method is the modified q-sample uniform scores test

99    (Stürzebecher et al., 1996; Stürzebecher et al, 1999). The modified q-sample uniform scores test is

100   applied to the ranks of the phases and amplitudes of multiple spectral bands, and has outperformed

101   various alternative methods when detecting the ABR. These include the original q-sample uniform

102   scores test and the q-sample analogue to Watson's U2 statistic (Stürzebecher et al, 1999), along with

103   the F for a Single Point (Fsp), Friedman's test, and Cochran's Q test (Cebulla et al., 2000). Moreover,

104   Cebulla et al (2006) have proposed various additional modifications to the q-sample uniform scores

105   test. These modifications have shown a high performance when detecting auditory steady state

106   responses (ASSRs), but have not yet been compared for ABR detection. The present paper will

107   investigate if these proposed modifications are also suitable for ABR detection.

108

109   Another promising multi-band detection method is the Hotelling's $T^2$ test (Hotelling, 1931), which has

110   outperformed both the Standard Deviation Ratio and the correlation coefficient (between two

111   replicates of the coherent average) when detecting ABRs in subject recorded data (Valdes et al.,

112   1987), along with the F for multiple points statistic (Fmp) when detecting ABRs extracted from quasi

113   ASSRs (Lachowska et al, 2012). The Hotelling's $T^2$ test has recently also been applied in the time

114   domain for detecting the slow cortical response (Golding et al., 2009; Carter et al., 2010; Chang et al.,

115   2012; Van Dun et al., 2012; Van Dun et al., 2015). These time-domain features (see section 3.1) have

116   not yet been evaluated for ABR detection, and may be a preferable alternative to frequency domain

117   analysis due to the broadband spectral content of the ABR.

118

119   Additional time-domain techniques for ABR detection methods that are of interest include the Fsp and

120   the Fmp, both of which can be tested for significance using F-distributions with $v_1$ and $v_2$ degrees of

121   freedom. A recurring complication, however, is that the degrees of freedom are typically unknown for

122   EEG data, and hence have to be assumed before statistical inference can be realized. Because false-

123   negatives (i.e. failure to detect an ABR response that is in fact present) are typically less detrimental

124   to the performance of ABR applications (ABR hearing screening tests in particular) than false-

125   positives, Elberling & Don (1984) have recommended a conservative approach (fewer false-positives

126    than the nominal target) by setting $v_1$ to 5. The drawback is a decrease in test power, which may result

127    in an increased cost of service delivery due to prolonged test times and/or increased false negative

128    rates.

129

130    An alternative approach for evaluating the significance of the Fsp (and the Fmp) has been proposed

131    by Lv et al (2007). Lv et al use bootstrapping (Efron & Tibshirani, 1993 - see also section 3.4) to

132    approximate the statistic's underlying null distribution. The approximated null distribution can then be

133    used for statistical inference without needing to explicitly estimate or assume the degrees of freedom

134    of the data. It is therefore hypothesized that evaluating statistical significance with the bootstrap

135    approach, as opposed to using F-distributions with assumed degrees of freedom, would provide more

136    consistent results across datasets with degrees of freedom that may vary between individuals and

137    recordings.

138

139    The goal for this paper is to evaluate and compare the performance of various objective ABR

140    detection methods in terms of specificity, sensitivity and detection time by using simulations and a

141    small sample of normal-hearing adults. The methods selected for the analysis include two versions of

142    the modified q-sample uniform scores test, which use either the ranks or the actual values of the

143    phases and amplitudes of the Fourier components of multiple spectral bands (see section 3.2 for

144    details), Hotelling's $T^2$ test applied in either the time or the frequency domain (section 3.1), and both

145    the Fsp and the Fmp (section 3.3), which were evaluated using either F-distributions with assumed

146    degrees of freedom or with the bootstrap approach.

147

148

## 149    2. ABR and no-stimulus EEG data

150    The data used throughout this study consists of (i) a small sample of normal-hearing adults where

151    physiological hearing thresholds were estimated using click-evoked ABRs of various intensity levels,

152    thus yielding a wide range of ABR waveform morphologies, and (ii) a relatively large database of no-

153    stimulus EEG recordings. The database of no-stimulus EEG recordings was initially used to assess the

154    specificities of the methods (section 3.5), after which it was used in combination with the subject

155    recorded ABR data in simulations to assess sensitivity (section 3.6). The simulated data provides a

156    test-bed in which large amounts of well controlled data are generated to assess performance when

157    signal characteristics are repeatable. The next step was to assess sensitivities and detection times of

158    the methods using just the subject recorded ABR data (section 3.7), which reflect real-world features

159    of routine recordings. Extended clinical studies including data from participants with a range of

160    hearing impairments are beyond the scope of the current work, but should follow in progressing this

161    research further.

162

163

### 2.1 Subject recorded ABR data

165    The subject recorded ABR data, previously described in Lv et al (2007), was collected from 12

166    subjects (6 female and 6 males) ranging from 18 to 30 years of age. The stimulus was a rectangular

167    100 μs click delivered at a stimulus rate of 33.3 Hz through ER-2 insert phones (Etymotic, USA). The

168    click intensities ranged from 0 to 50 dB SL (sensation level, i.e. relative to individual hearing

169    thresholds) in steps of 10 dB. The behavioural thresholds were estimated using a simple 'up-down'

170    approach where the click intensity was reduced in steps of 10 dB for every correct response, and

171    increased in steps of 5 dB for every missed response. ABRs were recorded with the active electrode

172    placed at vertex, a reference electrode at the nape of the neck, and a ground electrode placed at mid-

173    forehead. Measurements were obtained at a sampling rate of 10 kHz using a Cambridge Electronic

174    Design (CED) micro 1401 data acquisition unit along with a CED 1902 amplifier. Electrode

175    impedances remained below 5 kΩ throughout the recording. The recordings were band-pass filtered

176    offline from 30 to 1500 Hz with a 3[rd]-order Butterworth filter. Each recording was furthermore

177    downsampled to 5 kHz, and an artefact rejection method was applied by discarding 15% of the

178    noisiest epochs, as determined by their mean square values. Approximately 3600 clicks were

179    delivered per subject and per stimulus condition, resulting in a minimum of 3000 epochs after artefact

180  rejection. The 30.03 ms intervals following the onset of each stimulus (henceforth referred to as

181  epochs) were saved for offline analysis.

182

183

184  ### *2.2 No-stimulus EEG recordings*

185  Recordings of spontaneous EEG background activity (no stimulus was used) were previously

186  collected by Madsen et al (2017) and Madsen (2010) from 17 subjects (12 male and 5 female) under

187  four conditions. The conditions were (i) *asleep*, where the subjects were asked to try and fall asleep,

188  though sleep was not confirmed, (ii) *still*, where the subjects were instructed to lie still with their eyes

189  closed, but not to fall asleep, (iii) *blink*, where the subjects were instructed to blink every 1 to 3

190  seconds as a circle appeared on a screen in front of them, and (iv) *move*, where the subjects were

191  asked to move according to a random animation, also shown on a screen in front of them.

192  Measurements were then obtained using a Compumedics Neuroscan II EEG amplifier at a sampling

193  rate of 20 kHz with three silver-silver chloride (Ag/AgCl) electrodes placed on the left mastoid, the

194  right cheek (ground), and the upper forehead (reference). The electrode impedances remained below 1

195  $k\Omega$ throughout the recording for all subjects.

196

197  In the present study, the background EEG recordings were band-pass filtered with a 3$^{rd}$-order

198  Butterworth filter from 30 to 1500 Hz, after which they were downsampled to 5 kHz. Each recording

199  was then structured into 30.03 ms epochs, and artefact rejection was applied by discarding 15% of the

200  noisiest epochs, as determined by their mean square values. A total of 149 continuous EEG recordings

201  were available, with an average of 6800 pre-processed epochs per recording, resulting in

202  approximately 8 hours' of EEG.

203

204  # 3. Methods

205  This section first provides a description of the ABR detection methods and the bootstrap approach,

206  after which the adopted methodologies for evaluating the specificity, sensitivity, and detection time of

207  the methods are described.

208

209  The data to which the methods are applied consists of ensembles of epochs (for details on how these

210  ensembles were pre-processed and constructed, the reader is referred to sections 2.1, 2.2, 3.5 and 3.6).

211  Each ensemble is structured according to matrix **D**:

212  *Equation 1.*

213  $$\mathbf{D} = \begin{bmatrix} d_{11} & \cdots & d_{1K} \\ \vdots & \ddots & \vdots \\ d_{N1} & \cdots & d_{NK} \end{bmatrix}$$

214  where N is the ensemble size, K is the number of samples per epoch, and $d_{ij}$ is the $j^{th}$ sample of the $i^{th}$

215  epoch. The mean epoch $\overline{E}$ (also known as the coherent average) is found by taking the K averages

216  across the columns. The frequency domain representation of **D** is furthermore obtained by taking the

217  Fast Fourier Transform (FFT) of each row. Features can then be extracted from either the time or

218  frequency domain representations of the data. Extracting L features from each epoch results in the

219  NxL-dimensional feature matrix **V**:

220  *Equation 2.*

221  $$\mathbf{V} = \begin{bmatrix} v_{11} & \cdots & v_{1L} \\ \vdots & \ddots & \vdots \\ v_{N1} & \cdots & v_{NL} \end{bmatrix}$$

222  where $v_{ij}$ is the $j^{th}$ feature extracted from the $i^{th}$ epoch.

223

## 3.1 The one-sample Hotelling's $T^2$ test

224

225  The one-sample Hotelling's $T^2$ test is the multivariate extension to Students t-test, and can be used to

226  test whether the means of L features are significantly different from L hypothesized values. In the

227  present work it is assumed that the expected values of the features are zero. The statistic itself is a

228  weighted sum of the L feature means where the weights are determined by the variances and

229  covariances of the features. These weights have the convenient property of normalizing the L means,

230  which allows features with different scales and units to be combined appropriately. The $T^2$ statistic is

231  given by (Rencher, 2001):

232  *Equation 3.*
$$T^2 = N\,(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)\mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^H$$

233  where $\bar{\mathbf{x}}$ is the L-dimensional vector of means (found by taking the means down the L columns of

234  **V)**, $\boldsymbol{\mu_0}$ is the L-dimensional vector of hypothesized values to test against, $\mathbf{S}^{-1}$ is the inverse of the

235  covariance matrix of the NxL-dimensional feature matrix $\mathbf{V}$, and the $^H$ superscript denotes the

236  Hermitian transpose. The $T^2$ statistic can then be transformed into an F statistic with $v_1$ and $v_2$ degrees

237  of freedom using:

238  *Equation 4.*
$$F = \frac{N - L}{L(N - 1)}T^2 \quad \sim F_{L,N-L}$$

239  where $v_1 = L$ and $v_2 = $ N-L. The significance of F can be determined with an F-distribution look-up

240  table, or by finding the area under an F-distribution with L and N-L degrees of freedom on the interval

241  0 to F.  Note that in order to calculate $\mathbf{S}^{-1}$, the number of epochs N should be larger than the number of

242  features L. Note also that when the features are highly correlated, that $\mathbf{S}^{-1}$ can be close to singular, in

243  which case rounding errors might occur. A solution would then be to use the pseudoinverse (e.g. the

244  Moore-Penrose pseudoinverse; Moore, 1920; Penrose, 1955) instead of the regular inverse.

245

*Time domain features*

When applied in the time domain, the features for the one-sample Hotelling's $T^2$ test consist of `time-voltage means' (TVMs), which are defined as mean voltages, calculated across short time-intervals within each epoch (see e.g. Golding et al., 2009; Carter et al., 2010; Chang et al., 2012; Van Dun et al., 2012; Van Dun et al., 2015). As an example, when NxL TVMs are extracted, then each epoch is divided into L segments of approximately equal duration, and the mean is taken across each segment, resulting in the NxL-dimensional feature matrix **V**. The length of each segment requires a compromise such that the segments are neither too long, thus covering both peaks and troughs (with an average value of approximately zero) nor too short, thus leading to poor statistical robustness and reduced sensitivity. Because the direct current component is removed from the EEG recordings with a high-pass filter, the expected values for the TVMs will be zero. The hypothesized values to test against (defined above as $\mu_0$) are therefore given as an L-dimensional vector of zeros.

*Frequency domain features*

When using the frequency domain approach, the features are the real and imaginary parts of the Fourier components of Q spectral bands (resulting in an Nx2Q-dimensional feature matrix **V**). Because the phases of each spectral band are expected to be uniformly distributed between 0 and $2\pi$ when no response is present, the expected values for the real and imaginary parts of each spectral band are again 0. The hypothesized values to test against are therefore given as a 2Q-dimensional vector of zeros.

### 3.2 The modified q-sample uniform scores test

The original q-sample uniform scores test (Mardia, 1972) is a non-parametric test that uses the ranks of the phases of the Fourier components of Q spectral bands to test whether the phases share the same

271    distribution. The modification proposed by Stürzebecher et al (1999) uses the ranks of the amplitudes

272    in addition to the ranks of the phases, and is given by:

273    *Equation 5.*

$$W^* = C \sum_{j=1}^{G} \left[ \left( \sum_{i=1}^{N} r_{ij} \cos \beta_{ij} \right)^2 + \left( \sum_{i=1}^{N} r_{ij} \sin \beta_{ij} \right)^2 \right]$$

274    where $r^{ij}$ is the rank of the amplitude of the $i^{th}$ Fourier component (obtained from the $i^{th}$ epoch) of the

275    $j^{th}$ spectral band. C is furthermore an additional scaling factor given by

276    *Equation 6.*

$$C = \frac{4}{Q^2(Q+1)^2} \frac{2}{N}$$

277    and $\beta_{ij}$ is given by:

278    *Equation 7.*

$$\beta_{ij} = \frac{a_{ij} 2\pi}{NQ}$$

279    where $a_{ij}$ is the rank of the phase of the $i^{th}$ Fourier component (obtained from the $i^{th}$ epoch) of the $j^{th}$

280    spectral band. This modification will henceforth be referred to as 'Modified q-sample (ranks)' (using

281    the same notation as Cebulla et al, 2006).

282

283    In addition to the Modified q-sample V2 test, the 'Modified q-sample V4' test (Cebulla et al., 2006) is

284    also included in the analysis. The latter uses the actual values of the phases and amplitudes as opposed

285    to their ranks, in which case $r_{ij}$ in equation 5 refers to the amplitude of the $i^{th}$ Fourier component of the

286    $j^{th}$ spectral band, and $\beta_{ij}$ to the (untransformed) phase value of the $i^{th}$ Fourier component of the $j^{th}$

287    spectral band. The significance of these statistics can furthermore be evaluated with pre-determined

288    critical values based on simulations (Stürzebecher et al, 1999; Cebulla et al, 2000; Cebulla et al,

289    2006). Deviating from the literature, the significance of the Modified q-sample V2 and V4 statistics in

290    this study are evaluated using the bootstrap, as opposed to using pre-determined thresholds generated

291    from no-stimulus data. How the critical decision thresholds might differ between the two approaches

292    is further considered in the discussion.

293

294

295    **_3.3 The Fsp and the Fmp_**

296    The Fsp and the Fmp are defined as the ratio between the variance of the mean epoch $\bar{E}$ (found by

297    taking the K averages across the columns of data matrix **D**) and the estimated variance of the EEG

298    background noise. For the Fsp, the variance of the EEG background noise is estimated by the 'single

299    point' (SP) variance, which is defined as the variance down a single column of data matrix **D**. The Fsp

300    is given by (Elberling & Don, 1984):

301    *Equation 8.*
$$Fsp = N \frac{VAR(\bar{E})}{VAR(SP)}$$

302    where VAR denotes sample variance and SP refers to the values along an arbitrarily chosen column of

303    **D**. For the Fmp, the variance of the EEG background noise is estimated by taking the average of

304    multiple 'SP variances' (the average of the variances of multiple columns of **D**). The Fmp is given by

305    (Martin et al., 1994):

306    *Equation 9.*
$$Fmp = N \frac{VAR(\bar{E})}{\frac{1}{M} \sum_{j=1}^{M} VAR(SP_j)}$$

307    where $VAR(SP_j)$ is the variance of the $j^{th}$ included column of **D**, and M is the number of columns (of

308    **D**) to include.

309

310    Under the null hypothesis of no response present, it is assumed that the Fsp and the Fmp follow F-

311    distributions with $v_1$ and N-1 degrees of freedom. The latter is justified by assuming that epochs are

312    sufficiently distant in time to be uncorrelated (and thus independent for normally distributed data).

313    When the spectrum of the coherent average $\overline{E}$ is white, then the K samples within the coherent

314    average can also be considered independent and $v_1$ will equal K. The finite frequency content of EEG

315    background activity, however, introduces correlations between the samples, which makes the true

316    degrees of freedom difficult to estimate. A conservative recommendation, i.e. a FPR smaller than the

317    nominal alpha level of the test, for $v_1$ is 5 (Elberling and Don, 1984). As an alternative, the Fsp and

318    Fmp can be evaluated with the bootstrap approach.

319

320

### 3.4 Bootstrapping

322    Bootstrapping is a resampling with replacement procedure that can be used to construct a reference

323    distribution so that statistical inference can be performed (Efron & Tibshirani, 1993). For evoked

324    response detection, the goal is to construct the null distribution of some statistic by repeatedly

325    drawing ensembles of epochs from the continuous EEG record, and calculating the statistic of interest

326    on each new ensemble. Each bootstrapped ensemble should therefore represent the no-response

327    condition, achieved by randomly selecting N segments of K samples from within the continuous EEG

328    without regard to where the stimuli occur (Lv et al., 2007). Note that the selected segments may

329    overlap, in accordance with the principles of bootstrapping where samples are picked at random with

330    replacement, i.e. without removing that data from what can be picked later. The null distribution is

331    then approximated by calculating the statistic in question from many bootstrapped ensembles (1000

332    ensembles were used for this study). Finally, the statistic is also calculated from the original ensemble

333    of epochs, and its significance is evaluated by finding its location (percentile) along the bootstrapped

334    null distribution.

335

336

### 3.5 Specificity assessment

338    A methods FPR, or 1-specificity, is defined in the current work as the percentage of significant test

339    outcomes when no response is present (note again that this definition differs from studies that aim to

340     detect a clinical disorder, e.g. hearing loss, where sensitivity commonly refers to the detection of the

341     disorder). The FPRs of the methods were evaluated for different ensemble sizes using the pre-

342     processed recordings of EEG background noise (no stimulus was used) described in section 2.2. The

343     ensemble sizes selected for the analysis were 50, 100, 175, 275, 375, 500, 650, and 800 epochs, which

344     were chosen based on results from the sensitivity assessment (see section 3.6). For each ensemble

345     size, the EEG recordings were decomposed into ensembles of (consecutive) 30.03 ms epochs,

346     resulting in a total of 20197, 10060, 5717, 3606, 2640, 1967, 1500, and 1187 ensembles with

347     ensemble sizes of 50, 100, 175, 275, 375, 500, 650, and 800 respectively. Note that no further

348     distinction was made between EEG recordings obtained under different noise conditions. The latter

349     keeps the results concise, and is justified by realizing that all four conditions occur in clinical practice,

350     and that, ideally, the methods should perform adequately under each of them. The detection methods

351     were then applied to the initial 15 ms windows of the 30.03 ms epochs within each ensemble. For the

352     frequency domain methods, the spectral resolution was first increased to 40 Hz by extending each 15

353     ms segment to 25 ms with zero-padding. The frequency domain methods were then applied to all

354     spectral bands between 80 and 600 Hz. For the Modified q-sample V2 and V4 tests, averaging was

355     used (prior to calculating the FFT) to compress each ensemble into blocks of sub-averages, as

356     recommended by Cebulla et al (2000) (As noted by one of the reviewers, it is worth emphasizing that

357     these recommendations were formulated for ABR detection, and that in later publications on ASSR

358     detection, averaging is not advocated, see Cebulla et al, 2006). In this study, averaging was performed

359     across blocks of 25 epochs so that no epochs were excluded from the analysis (each ensemble size is a

360     multiple of 25), which hence compressed each ensemble into $\frac{N}{25}$ sub-averages. With respect to the

361     time domain methods, a total of 25 TVMs were used for the Hotelling's $T^2$ test, which were spread

362     equally across the 15 ms analysis window. The choice for 25 TVMs was based on additional

363     simulations, which showed a robust performance for the Hotelling's $T^2$ test when using anything

364     between ~20 and ~40 TVMs. These simulations were similar to the ones described in section 3.6

365     below, but used an alternative set of ABR templates for simulating a response (obtained from the

366     coherent averages of the subject data described in Elberling et al, 2010). The column index (of data

367    matrix **D**) for calculating the single point variance for the Fsp was furthermore arbitrarily set to 30,

368    and the number of columns to include in the Fmp was set to 75 (corresponding to the full analysis

369    window, or 15 ms). The significance of the Fsp and Fmp was evaluated using either F-distributions

370    with 5 and N-1 degrees of freedom (denoted by 'Fsp 5 dof' and 'Fmp 5 dof' respectively) or with the

371    bootstrap approach (denoted by 'Fsp bootstrapped' and 'Fmp bootstrapped' respectively).

372

373

374    *3.6 Sensitivity assessment using simulations*

375    A methods TPR, or sensitivity, is defined as the percentage of significant test outcomes (ABR

376    responses detected) when a response is present, which should of course be as high as possible for

377    some set FPR. In this study, sensitivity was assessed using both simulations and subject recorded

378    ABR data. Simulations were included as these allow a large number of tests to be performed, which

379    allows powerful comparisons to be drawn amongst the methods. This is important for the present

380    study as the analysis regarding the subject recorded ABR data (section 3.7) was based on just 12

381    subjects.

382

383    The data used for the simulations consists of (i) the pre-processed recordings of EEG background

384    noise (see section 2.2), along with (ii) the coherent averages from the subject recorded ABR data that

385    contained a clear response. The latter was determined through visual inspection by an experienced

386    audiologist. As guidance for determining the presence of a clear response, the audiologist inspected

387    the repeatability of the waveform by comparing two replicates of the coherent average (obtained by

388    averaging across epochs 1 to 1500, and again across epochs 1501 to 3000). The audiologist also used

389    the 3 to 1 signal to noise criterion as additional guidance (see Sutton et al, 2003), but was ultimately

390    left free to decide whether a response was present or not. This process resulted in a total of 34 ABR

391    templates with a clear response: 4, 7, 8, 7, and 8 from the 10, 20, 30, 40, and 50 dB SL conditions

392    respectively. Data was then assembled by randomly selecting N consecutive epochs from within a

393     randomly selected recording of EEG background noise, and adding a randomly selected and rescaled

394     ABR template to all epochs within the ensemble. The ensemble size N took values of 50, 100, 175,

395     275, 375, 500, 650 and 800 epochs, which were chosen based on the results from a pilot simulation

396     that showed a good coverage of TPRs across methods when using these values. The scaling factor was

397     furthermore chosen so that the signal to noise ratio (SNR) was -23 dB, which was calculated

398     according to:

399

400     *Equation 10.*

401     $\text{SNR} = 10\log_{10}\left(\dfrac{\text{P}_{\text{Signal}}}{\text{P}_{\text{Noise}}}\right)$

402     where $\text{P}_{\text{Signal}}$ is the mean square of the scaled ABR waveform, and $\text{P}_{\text{Noise}}$ the mean square of the

403     ensemble of N epochs (prior to adding the ABR waveform, and treated as a continuous recording).

404     The SNR of -23 dB was based on a brief analysis of the subject recorded ABR data, which showed

405     that the responses were in the proximity of -23 dB. The latter was similarly calculated with equation

406     10, with $\text{P}_{\text{Signal}}$ now being the mean square of the coherently averaged ABR (calculated across all 3000

407     epochs from the subject and dB SL condition in question), and $\text{P}_{\text{Noise}}$ the mean square of the epochs

408     when treated as a continuous recording. A total of 10 000 tests were performed for each ensemble size

409     using the same detection methods and features as those described in section 3.5.

410

411     ### *3.7 Sensitivity and detection time assessment using subject recorded ABR data*

412     The sensitivities and detection times of the methods were further evaluated using just the subject

413     recorded ABR data. The methods were applied to the initial 1-16 ms segments of the epochs (the first

414     ms was excluded to avoid potential contaminations from a stimulus artefact), which was repeated for

415     each subject and each stimulus condition. The methods were applied to the data sequentially, every 50

416     epochs, from 50 epochs onwards. To clarify - a test was first performed using an ensemble size of 50,

417     then again using an ensemble size of 100, etc., until all 3000 epochs had been analysed (a total of 60

418    tests, per subject, and per dB SL condition). The detection methods and features selected for the

419    analysis were the same as those described in section 3.5.

420

421

422    ## *4. Results*

423    This section presents the results regarding the specificity assessment (section 3.5), the sensitivity

424    assessment using simulations (section 3.6), and the sensitivity and detection time assessment using the

425    subject recorded ABR data (section 3.7).

426

427    *4.1 Specificity assessment*

428    The FPRs of the methods (using an alpha of 0.01) for the no-stimulus condition are presented in Table

429    1 for different ensemble sizes N. The upper and lower boundaries for significant deviations (**p<0.05**)

430    from the expected 1% FPRs were found using the binomial distribution (see appendix). Results show

431    that the FPRs of `Fsp 5 dof' and `Fmp 5 dof' were significantly (**p<0.05**) lower than 1%, as predicted

432    by Elberling & Don (1984). The remaining methods appear to show a slight tendency towards a

433    higher than expected FPR, which was significant (**p<0.05**)  for: 'T2 Time' (for N=100 and N=375) ,

434    'T2 Freq' (for N=375 and N=650), 'Fsp bootstrapped' (N=50, N=100, N=175, and N=375), 'Fmp

435    bootstrapped' (N=100 and N=175), and 'Modified q-sample V2' (N=375 and N=500). Although these

436    deviations are relatively small (all remained below 2%), a higher than expected FPR can be

437    worrisome for some ABR applications. Additional simulations were therefore performed to further

438    test and explore why the methods appear to show a higher than expected FPR. The data for these

439    simulations consisted of a large amount (50 000 recordings) of Gaussian-distributed coloured noise

440    with similar spectral content to real EEG background noise, and with stationary variance and a true

441    mean of zero. The resulting data was pre-processed and analysed as described in sections 2.2 and 3.5,

442    i.e. the settings were the same as those used for evaluating the real EEG background activity. The

443    resulting FPRs were in the range of 1.15% to 1.2% for an expected FPR of 1%. Although very small,

444     the deviations were significant (**P<0.01**) as the statistical power was high (50 000 tests were

445     performed). Note that all underlying statistical assumptions for these simulations were met, except

446     potentially the independence assumption between epochs, which suggests that the slightly higher than

447     expected FPR can be attributed to a violation of this assumption. The latter is further addressed in the

448     discussion.

449

450     - INSERT TABLE 1 -

451

452     *4.2 Sensitivity assessment using simulations*

453     The detection rates of the methods (using an alpha of 0.01) are presented in Figure 1 as a function of

454     the ensemble size N. The best performances (highest TPR) is noted consistently for the HT2 tests,

455     followed by the modified Q-sample tests, the bootstrapped Fmp and Fsp, and lastly by the Fsp and

456     Fmp evaluated with F-distributions with assumed degrees of freedom. As the latter have a lower FPR

457     also (see Table 1), a reduced TPR also might be expected. Moreover, the Fmp and Fsp use only the

458     SNR (i.e. average power values) and it is thus not surprising that they are less sensitive. Note that a

459     potential danger of using detection rates for comparisons in sensitivity is that methods with higher

460     FPRs are given an unfair advantage over those that are more conservative and have a lower FPR. As

461     shown by Table 1, the FPRs were all close to the expected 1% (with the exception of 'Fsp 5 dof' and

462     'Fmp 5 dof '), which suggests that the comparison in sensitivity was fair. The latter was verified by

463     finding the critical alpha values under which the methods obtained a FPR of exactly 1%, and

464     replotting the resulting detection rates. Results (not presented) were almost identical to those

465     presented in Figure 1 (again, with the exception of 'Fsp 5 dof' and 'Fmp 5 dof '). With respect to 'Fsp

466     5 dof' and 'Fmp 5 dof', their sensitivity was greatly improved (to values similar to those seen with the

467     bootstrapped method) by using the adjusted critical alpha values. It also might be noted that the

468     differences in FPR shown in Table 1 do not consistently explain the differences in TPR between the

469     methods.

470

471     - INSERT FIGURE 1 –

472

473     *4.3 Sensitivity and detection time assessment using subject recorded*

474     *ABR data*

475     The detection rates for an ensemble size of 3000 are presented in Figure 2 for the 0, 10, 20, and 30 dB

476     SL conditions (the 40 and 50 dB SL conditions are excluded as all methods obtained a 100% detection

477     rate here). The required time for detecting a response was then found by finding the number of stimuli

478     (expressed in seconds) required for the p-value to drop and remain below the 0.01 threshold for the

479     remainder of the test. The additional requirement that the p-value remains below the 0.01 threshold

480     ensures that the FPR is not inflated due to multiple tests being performed. If a test did not drop below

481     the 0.01 significance threshold, then the full ~90 seconds test time was used (corresponding to 3000

482     epochs), which may have resulted in an underestimation of the required test time in the case of a false

483     negative. The mean of the resulting detection times (taken across subjects) are presented in Figure 3

484     as bar graphs for each method, and dB SL condition. HT2 consistently showed the best performance.

485

486     - INSERT FIGURE 2 -

487

488     - INSERT FIGURE 3 -

489

490     Visually inspecting the distributions of the detection rates and detection times showed that both were

491     strongly non-Gaussian, which was confirmed with the Kolmogorov-Smirnov goodness of fit test

492     (**p<0.01** for all distributions). Non-parametric statistical analysis was therefore used to test whether

493     the discrepancy amongst the methods in terms of detection rates and detection times was significant.

494     With respect to detection rates, Cochran's Q test was first used to test for equivalence in performance

495     across all 8 methods for each dB SL condition. Results indicate a significant difference in

496     performance for the 10 (**p < 0.01**) and 20 (**p < 0.05**) dB SL conditions. As a follow-up, Fishers exact

497     test was used to draw pairwise comparisons amongst the methods for the 10 and 20 dB SL conditions.

498   Results show that the performance of T2 Time and T2 Freq both differed significantly (**p<0.05**) from

499   'Fmp 5 dof' and 'Modified q-sample V4' for the 10 dB SL condition. The latter is presented in Figure

500   2 using asterisks and crosses, where an asterisk denotes a significant discrepancy with T2 Time, and a

501   cross a significant discrepancy with T2 Freq. Comparisons between the remaining methods were not

502   significant. Similarly, with respect to detection times (Figure 3), non-parametric statistical analysis

503   was first used to test for equivalence in performance across all 8 methods (now using Friedman's

504   test), per dB SL condition. Results indicate a significant difference in performance for the 10, 20, 30,

505   40, and 50 dB SL conditions (all **p < 0.001**). The Wilcoxon rank sum test was then used to draw

506   pairwise comparisons between all methods, for the 10, 20, 30, 40, and 50 dB SL conditions. The

507   majority of the comparisons were again not significant, with the exception of the Hotelling's T2 test.

508   Significant (**p<0.05**) discrepancies between T2 Time, T2 Freq, and the remaining methods are again

509   represented by asterisks and crosses in Figure 3, with asterisks denoting a significant discrepancy with

510   T2 Time and crosses with T2 Freq.

511

# 5. Discussion

513   This study used simulations and subject recorded data to compare the specificity, sensitivity and

514   detection time of various objective ABR detection methods. With respect to specificity, although the

515   FPRs mostly fell within the lower and upper 95% confidence intervals of the expected FPR, a slightly

516   higher than expected FPR was observed. The Fsp and Fmp evaluated with theoretical F-distributions

517   were the exception, both of which showed consistently lower than expected FPRs. In terms of

518   sensitivity and detection time, the Hotelling's $T^2$ test came out on top in both the simulations and the

519   subject recorded data. The results regarding these properties (specificity, sensitivity, and detection

520   time) are now discussed in more detail.

521

522

## 5.1 Specificity

The results from the specificity assessment (Table 1) suggest a slightly higher than expected FPR for most methods (excluding 'Fsp 5 dof' and 'Fmp 5 dof'). Although the deviations were small, this is a concern for many ABR applications where a higher than expected FPR can potentially have severe repercussions. In ABR hearing screening applications, for example, a higher than expected FPR can result in additional cases of undetected hearing loss (ABR responses are falsely detected), which (when left untreated) have been associated with an impaired language development in children (Ramkalawan & Davis, 1991), along with various other more obvious handicaps, such as discrimination, less effective education, a reduced life expectancy, and higher unemployment rates (Miziara, 2012), to name a few. Supplementary simulations (see section 4.1) were therefore performed to explore why the FPRs appear to be higher than expected. Results suggest that the increased FPR can be attributed to a violation of the independence assumption between epochs. For recorded EEG data, which is known to be dominated by low frequencies with an approximate $\frac{1}{f^\alpha}$ spectrum (with $\alpha \approx 1$, see Pritchard, 1992), a similar violation can be expected, with the extent of the violation depending primarily on the cut-off frequency of the high-pass filter and the inter-epoch interval (i.e. the stimulus rate). It can thus be expected that increasing the high-pass cut-off frequency or decreasing the stimulus rate would reduce the long term correlations between epochs, and increase their independence. Results from an additional post-hoc analysis indeed show no significant (**p<0.05)** deviations from the expected 1% FPR when repeating the specificity assessment (section 3.5) with an adjusted high-pass cut-off frequency of 100 Hz. An alternative solution to increased FPRs may be to adjust the significance level (the alpha value) of the test. Post hoc analysis showed that FPRs of exactly 1% (across all ensemble sizes in Table 1) could be obtained when using the following alpha values: 0.0087 (T2 Time), 0.0088 (T2 Freq), 0.021 (Fsp 5 dof), 0.0321 (Fmp 5 dof), 0.0071 (Fsp bootstrapped), 0.0074 (Fmp bootstrapped), 0.0088 (Modified q-sample V2), and 0.009 (Modified q-sample V4). Although these adjusted alpha values would result in a small loss of sensitivity, they may be a safer option when detecting ABRs in practice.

550 It should furthermore be emphasized that the adjusted alpha values for 'Fsp 5 dof' and 'Fmp 5 dof'

551 are expected to be more susceptible to the high-pass cut-off frequency than the remaining methods

552 due to an additional violation of the independence assumption amongst samples within epochs. In

553 particular, increasing the high-pass cut-off frequency results in fewer low frequency components,

554 which reduces the correlations amongst the samples within the epochs, thus increasing the degrees of

555 freedom of the data and removing it farther from the assumed degrees of freedom $v_1=5$. It can

556 therefore be expected that the performance of 'Fsp 5 dof' and 'Fmp 5 dof' would be even more

557 conservative when the high-pass cut-off frequency is increased, which was confirmed when repeating

558 the specificity assessment (section 3.5) with an adjusted high-pass cut-off frequency of 100 Hz. In

559 particular, the degrees of freedom $v_1$ for the no-stimulus EEG recordings (section 2.2) ranged from 3

560 to 15 (with mean 8.4 and standard deviation 2.6) when using a high-pass cut-off frequency of 30 Hz,

561 and from 3 to 20 (with mean 11.3 and standard deviation 4.9) when using a high-pass cut-off

562 frequency of 100 Hz (the latter was achieved by fitting F-distributions to each bootstrapped null

563 distribution, and finding the best fitting function). It might be noted that the conservative estimate of 5

564 degrees of freedom was originally intended for the higher cut-off frequency of 100 Hz (Elberling &

565 Don, 1984). Note also that the Hotelling's $T^2$ test and the bootstrapped statistics are immune to

566 independence violations amongst samples within epochs (but not between epochs). In particular, the

567 Hotelling's $T^2$ test accounts for correlated samples within epochs by scaling the features by their

568 covariance matrix (see methods section), whereas the bootstrapped statistics account for correlated

569 samples by resampling on an epoch to epoch basis, which preserves the correlations between samples

570 within epochs. This further allows the bootstrapped confidence intervals to more accurately reflect

571 test-dependent factors, such as the EEG background noise, the electrode impedances, and ultimately

572 the degrees of freedom of the data. The latter is important for many ABR applications where the

573 objective detection methods are expected to perform adequately across EEG recordings with varying

574 degrees of freedom. It is hence hypothesized that the Hotelling's $T^2$ test and the bootstrapped statistics

575 would provide more consistent results relative to 'Fsp 5 dof' and 'Fmp 5 dof' across a wider range of

576 test conditions. A similar argument might be made in favour of the bootstrap approach over the use of

577 pre-determined thresholds generated from no-stimulus data (see e.g. Stürzebecher et al, 1999; Cebulla

578   et al, 2000; Cebulla et al, 2006), i.e. pre-determined thresholds may not generalize well across data

579   sets and test conditions, whereas the bootstrap approach estimates confidence intervals specifically for

580   the recording in question.

581

582

583   **5.2 Sensitivity and detection time**

584   With respect to sensitivity (the proportion of correctly identified responses) and detection time, these

585   should ideally be as high and low as possible respectively for some set FPR. In ABR audiometry, for

586   example, one would expect thresholds to decrease as the sensitivity of the detection method is

587   increased, which may lead to greater convergence between behavioural and estimated hearing

588   thresholds. In terms of reduced test time, one would expect an increased sensitivity to result in (i) a

589   decreased cost of service delivery, (ii) reduced patient discomfort, and (iii) a smaller time window

590   within which noise artefacts can be introduced to the data. Reduced detection times would be

591   particularly beneficial in patients who cannot cooperate, such as infants or some with dementia.

592

593   In this study, sensitivity was evaluated using detection rates, which have the desirable properties of

594   being intuitive and simple. However, as noted in results section 4.2, a potential risk of using detection

595   rates is that methods with higher FPRs receive an unfair advantage over those with lower FPRs (the

596   latter is most notably the case for 'Fsp 5 dof' and 'Fmp 5 dof', which were indeed designed to have

597   lower FPRs). The problem can be resolved by adjusting the nominal alpha values so that the FPRs are

598   equal across methods. Note however that although this allows for a more fair comparison, it is not

599   necessarily a realistic one as adjustment of the FPR may need to be carried out on an individual basis

600   using prior knowledge that is not always available. Results from the simulations nevertheless suggest

601   an advantage for the Hotelling's $T^2$ test when using both the adjusted and unadjusted critical alpha

602   values (results section 4.2). This is further supported by FPRs in Table 1: consistent differences in

603   detection rates (Figure 1) cannot be readily explained from relatively inconsistent FPRs. Results from

604   the subject recorded ABR data similarly suggest an overall advantage for the Hotelling's $T^2$ test

605 (Figures 2 and 3). The relative discrepancy in performance amongst the remaining methods, however,

606 is less clear, which can likely be attributed to a small sample size of just 12 subjects.

607

608

609 With respect to the frequency domain features for the Hotelling's $T^2$ test, it is worth noting that these

610 are essentially the same as those used by the Modified q-sample V4 test (the Hotelling's $T^2$ test is

611 applied to the real and imaginary parts of the Fourier components, whereas the Modified q-sample V4

612 test is applied to the phases and amplitudes), and yet a relatively large discrepancy in performance

613 was still observed. This can likely be attributed to the way in which features are weighted and

614 combined into a single statistic. In particular, the Hotelling's $T^2$ test weights the features according to

615 their variance and covariance, whereas the Modified q-sample V4 test does not. The latter results in

616 an L-dimensional hyper-ellipsoid (centred at features means $\bar{x}$) as H0 rejection region for the

617 Hotelling's $T^2$ statistic, where the shape of the ellipsoid is determined by the variance and covariance

618 of the features. Having an ellipsoid as rejection region means that the null hypothesis is more easily

619 rejected in some directions relative to others, meaning it has the potential of providing a more

620 powerful test relative to, for example, a spherical rejection region.

621

622 Based on the preceding paragraph, an identical performance between the Modified q-sample V4 test

623 and the Hotelling's $T^2$ test might be expected when applied to uncorrelated features with equal

624 variance, which was tested with additional simulations. In particular, simulations described in

625 Stürzebecher et al (1999) and Cebulla et al (2000) were implemented, which used Gaussian zero mean

626 white noise with stationary variance to represent the EEG background noise, along with a sinewave

627 multiplied with a Gaussian window for representing a response. The detection methods included for

628 these simulations were (i) the original q-sample uniform scores test (Mardia, 1972), (ii) both the

629 Modified q-sample V2 and V4 tests (Stürzebecher et al, 1999; Cebulla et al, 2006), and (iii) the

630 Hotelling's $T^2$ test using the frequency domain approach. As predicted, the Hotelling's $T^2$ test and the

631 Modified q-sample V4 test both came out on top in terms of sensitivity (with very similar

632 performances), followed by the Modified q-sample V2 test (using ranks rather than measured values),

633 and lastly by the original q-sample uniform scores test (which only uses phase ranks).

634

635

636 **Study limitations**

637 In this study, the investigators strived to present a comprehensive and fair comparison between

638 various time and frequency domain ABR detection methods in terms of their sensitivity, specificity,

639 and detection time. Whenever possible, feature selection and pre-processing parameters were based

640 on recommendations or findings from the literature. That said, it remains to be seen whether the

641 results presented in this study generalize across alternative feature and data sets. Various additional

642 parameters worth investigating include the time window selected for the analysis, how many and

643 which spectral bands to include for the frequency domain methods, and the selection of TVMs for T2

644 Time. With respect to the latter, a total of 25 TVMs spread equally across a 15 ms analysis window

645 were used for this study. The choice for 25 TVMs was based on results from additional simulations

646 described in section 3.5, which showed a good performance for the Hotelling's $T^2$ test when using

647 anything between ~20 and ~40 TVMs. It is however worth noting that these simulations did not

648 distinguish between stimuli of different intensities. Hence, although 25 TVMs may be a relatively

649 robust set of features for ABR detection, it is not necessarily optimal, and it may be beneficial to use

650 more specific arrangements of TVMs depending on the type of stimulus and/or stimulus parameters

651 being used. A general rule of thumb is that the optimal number of TVMs will tend to increase along

652 with the number of peaks in the ABR, since consecutive time-domain peaks within the ABR would

653 cancel each other out when the number of TVMs is too low (Golding, 2009). Hence, when the

654 stimulus intensity is decreased, and ABR waves I and III begin to disappear (Hall, 2006), it may be

655 beneficial to use fewer TVMs.

656 # 5. Conclusion

657 Comparisons were drawn between various objective ABR detection methods in terms of specificity,

658 sensitivity, and detection time. Results from the specificity assessment suggest a tendency towards

659     slightly higher than expected FPR across methods, which likely can be attributed to a violation of the

660     independence assumption between epochs. With respect to sensitivity and detection time, the

661     Hotelling's $T^2$ test came out on top, which was primarily attributed to a more suitable weighting of the

662     features. Finally, bootstrapping was shown to improve the reliability of the Fsp and the Fmp, as

663     opposed to when test significance was evaluated using F-distributions with the recommended

664     assumption of 5 degrees of freedom.

665

672

673     **Declarations of interest**

674     The authors report no conflicts of interest.

675

676

677

678

679

680     # Appendix

681     **The binomial distribution**

682     A Bernoulli trial is a random experiment with exactly two possible outcomes, typically interpreted as

683     'success' and 'failure'. When X Bernoulli trials are performed and the probability of a successful trial

684     is P, then the binomial distribution gives the probability densities of observing x successful trials. The

685     distribution is given by:

$$B(x|X,P) = \frac{X!}{x!\,(X-x)!} P^x \, (1-P)^{X-x}$$

686     For the specificity assessment, the number of Bernoulli trials X was set to the number of ensembles

687     tested, and the probability of a successful trial P to the expected probability of observing a false

688     positive (P=0.01). The resulting distribution was used to find the 95% confidence intervals for the

689     expected number of false positives.

690

691

692

693

694

695

696

697

698

699

700     # References

701   Carter L., Golding M., Dillon H., Seymour J. 2010. The Detection of Infant Cortical Auditory Evoked

702   Potentials (CAEPs) Using Statistical and Visual Detection Techniques. *J Am Acad Audiol,* 21, 347-

703   356.

704

705   Cebulla M., Stürzebecher E., Wernecke K.D. 2000. Objective detection of auditory brainstem

706   potentials - Comparison of statistical tests in the time and frequency domain. *Scand Audiol*, 29, 44-51.

707

708   Cebulla M., Stürzebecher E., Elberling C. 2006. Objective detection of Auditory Steady State

709   Responses: Comparison of One-Sample and q-Sample Tests. *J Am Acad Audiol*, 17, 93-103.

710

711   Chang H.W., Dillon H., Carter L., van Dun B., Young S.T. 2012. The relationship between cortical

712   auditory evoked potential (CAEP) detection and estimated audibility in infants with sensorineural

713   hearing loss. *Int J Audiol*, 51, 663-670.

714

715   Efron B., Tibshirani R.J. 1993. An Introduction to the Bootstrap. Chapman & Hall/CRC. Editors: Cox

716   D.R., Hinkley D.V., Reid N., Rubin D.B., Silverman B.W. Boca Raton, London, New York,

717   Washington D.C.

718

719   Elberling C. 1976. Action Potentials Recorded from the Promontory and the Surface, Compared with

720   Recordings from the Ear Canal in Man. *Scand Audiol,* 5, 69-78

721

722   Elberling C. 1979. Auditory electrophysiology: spectral analysis of cochlear and brain stem evoked

723   potentials. A comment on: Kevanishvili and Aponchenko: "Frequency composition of brain stem

724   auditory evoked potentials". *Scand Audiol*, 8, 57-64.

725

726   Elberling C., Don M. 1984. Quality estimation of averaged auditory brainstem responses. *Scand*

727   *Audiol*, 13, 187-197.

728

729    Elberling C., Callø J., & Don M. 2010. Evaluating auditory brainstem responses to different chirp

730    stimuli at three levels of stimulation. J. Acoust. Soc. Am., 128 (1) 215-223.

731

732    Golding M., Dilon H., Seymour J., Carter L. 2009. The detection of adult cortical auditory evoked

733    potentials (CAEPs) using an automated statistic and visual detection. *Int J Audiol*, 48, 833-842.

734

735    Hall J.W., 2006. New Handbook of Auditory Evoked Responses. Pearson, 1[st] edition.

736

737    Hotelling H. 1931. The Generalization of Student's Ratio. *Ann Math Statist*, 2, 360-378.

738

739    Kevanishvili Z., Aphonchenko V. 1979. Frequency composition of the brain-stem auditory evoked

740    potential. *Scand Audiol*, 8, 51-55.

741

742    Lachowska M., Bohórquez J., Ozdamar O. 2012. Simultaneous acquisition of 80 Hz ASSRs and

743    ABRs from quasi ASSRs for threshold estimation. *Ear Hear*, 33, 660-671.

744

745    Lv J., Simpson D.M., Bell S.L. 2007. Objective detection of evoked potentials using a bootstrap

746    technique. *Med Eng Phys*, 29, 191-198.

747

748    Madsen S.M.K. 2010. Accuracy of averaged auditory evoked potential amplitude and latency

749    estimates. Msc thesis, Dept. Elec. Eng., Tech. Uni of Denmark.

750

751    Madsen S.M.K., Harte J.M., Elberling C., Dau T. 2017. Accuracy of averaged auditory brainstem

752    response amplitude and latency estimates. *Int. Jnl. Audiol,* 1-9.

753

754    Mardia K.V. 1972. Statistics of directional data. 1972. London: Academic Press.

755

756    Martin W.H., Schwegler J.W., Gleeson A. L., Shi Y.B. 1994. New techniques of hearing assessment.

757    *Otolaryngol Clin North Am*, 27, 487-510.

758

759    Miziara I.D., Miziara C.S., Tsuji R.K., Bento R.F. 2012. Bioethics and medical/legal considerations

760    on cochlear implants in children. *Braz. J. Otorhinolaryngol*, 78: 70-79.

761

762    Moore E.H. 1920. On the reciprocal of the general algebraic matrix. *Bulletin of the American*

763    *Mathematical Society,* 26, 394–395.

764

765    Penrose R. 1955. A generalized inverse for matrices, *Proceedings of the Cambridge Philosophical*

766    *Society,* 51, 406–413

767

768    Pritchard, W. (1992). The brain in fractal time: 1/f-like power spectrum scaling of the human

769    electroencephalogram. *Intern. J. Neurosci*, 66, 119–129.

770

771    Ramkalawan T.W., Davis A.C. 1992. The effects of hearing loss and age of intervention on some

772    language metrics in young hearing-impaired children. *Br. J. Audiol*. 26: 97-107.

773

774    Rencher A.C. 2001. Methods of Multivariate Analysis. Second Edition. John Wiley & Sons, Inc.

775    Chapter 5, page 118.

776

777    Stürzebecher E., Cebulla M., Baag M., Thie R. 1996. Verfahren zur automatischen

778    Hörschwellenbestimmung, insbesondere bei Neugeborenen und Kleinkindern. Patent application

779    PCT/DE96/02453.

780

781    Stürzebecher E., Cebulla M., Wernecke K. 1999. Objective response detection in the frequency

782    domain: comparison of several q-sample tests. *Audiol Neurootol*, 4, 2-11.

783

784    Sutton G., Lightfoot G., Stevens J., Booth R., Brennan S., Feirn R., Meredith R. Guidance for

785    Auditory Brainstem Response testing in babies. March 2013, Version 2.1.

786

787    Suzuki T., Sakabe N., Miyashita Y. 1982. Power Spectral Analysis of Auditory Brain Stem Responses

788    to Pure Tone Stimuli. *Scand Audiol*, 11, 25-30.

789

790    Valdes-Sosa M.J., Bobes M.A., Perez-Abalo M.C., Perera M., Carballo J.A., Valdes-Sosa P. 1987.

791    Comparison of Auditory-Evoked Potential Detection Methods Using Signal Detection Theory.

792    *Audiology*, 26: 166-178.

793

794    Van Dun B., Carter L., Dillon H. 2012. Sensitivity of cortical auditory evoked potential (CAEP)

795    detection for hearing-impaired infants in response to short speech sounds. *Audiol Res*, 2e13, 65-76.

796

797    Van Dun B., Dillon H., Seeto M. 2015. Estimating Hearing Thresholds in Hearing-Impaired Adults

798    through Objective Detection of Cortical Auditory Evoked Potentials. *J Am Acad Audiol*, 26, 370-383.

799
800
801
802
803
804
805
806
807
808
809
810    **Tables**

811     Table 1. The percentage of significant (**p<0.01**) tests for the no-response condition, per method and

812     per ensemble size. Significant deviations (**p<0.05**) from the expected 1% FPR are indicated by a red

813     asterisk.

814

815

816

817

818

819

820

821

822

823

824

825

826

827     **Figures**

828     Figure 1. The detection rates of the methods (using an alpha of 0.01) as a function of the ensemble

829     size when detecting a simulated -23 dB response (see section 3.6 for details).

830

831     Figure 2. The percentage of detected responses (**p<0.01**) for each method, presented as bar graphs for

832     each dB SL condition. Non-parametric statistical analysis was used to test whether the discrepancy

833     between methods was significant (see section 4.3 for details). The majority of the comparisons were

834     not significant, with the exception of the Hotellings $T^2$ test. Significant discrepancies (**p<0.05**) with

835     T2 Time are indicated by an asterisk (placed above the bar of the corresponding method), whereas

836     significant discrepancies with T2 Freq are indicated by a cross.

837

838     Figure 3. The mean of the detection times of the methods (calculated across subjects), presented as

839     bar graphs for each per dB SL condition. Non-parametric statistical analysis was used to test whether

840     the discrepancy between methods was significant (details presented in section 4.3). The majority of

841     the comparisons were not significant, with the exception of the Hotellings $T^2$ test. Significant

842     discrepancies (**p<0.05**) with T2 Time are indicated by an asterisk, placed above the bar of the

843     corresponding method. Significant discrepancies with T2 Freq are indicated by a cross, similarly

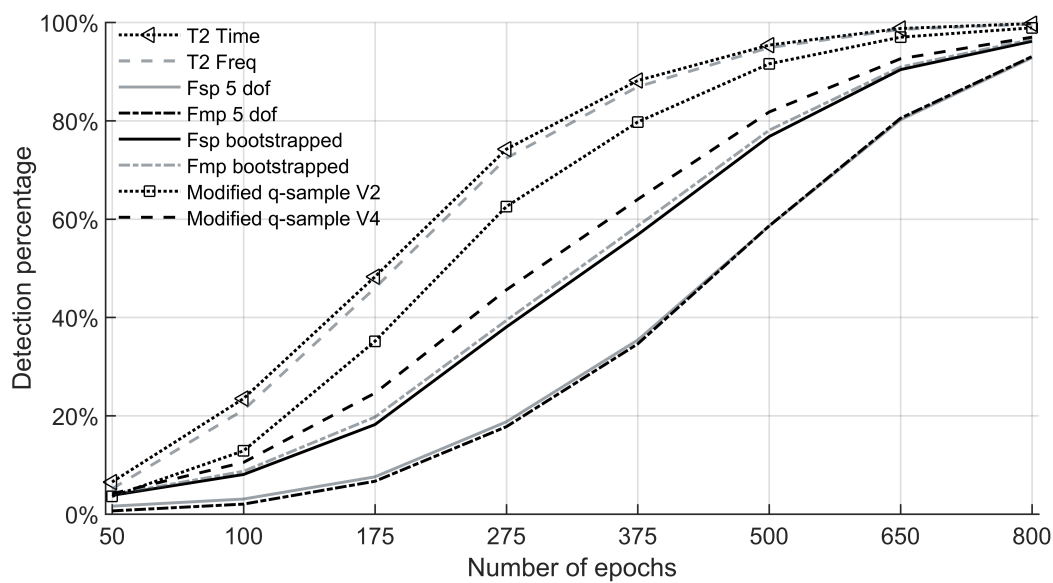844     placed above the bar of the corresponding method.

845

846

847

848

849

850

851

852     Figure 1

853

854

855

856

857

858

859

860

861

862

863

864

865

866     Figure 2

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881 Figure 3

Chesnaye et al: Objective measures for detecting the ABR: comparisons in specificity, sensitivity, and detection time

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897    Table 1

Chesnaye et al: Objective measures for detecting the ABR: comparisons in specificity, sensitivity, and detection time

| Number of epochs per ensemble -> | 50 | 100 | 175 | 275 | 375 | 500 | 650 | 800 |
|---|---|---|---|---|---|---|---|---|
| **T2 Time** | 1.08% | 1.25%* | 0.98% | 1.33% | 1.48%* | 0.92% | 1.13% | 1.26% |
| **T2 Freq** | 1.09% | 1.08% | 1.14% | 1.19% | 1.59%* | 1.32% | 1.73%* | 1.26% |
| **Fsp 5 dof** | 0.54%* | 0.53%* | 0.51%* | 0.5%* | 0.80% | 0.56%* | 0.4%* | 0.51%* |
| **Fmp 5 dof** | 0.23%* | 0.36%* | 0.37%* | 0.44%* | 0.61%* | 0.56%* | 0.33%* | 0.42%* |
| **Fsp bootstrapped** | 1.15%* | 1.27%* | 1.4%* | 0.94% | 1.44%* | 1.17% | 1.27% | 1.52% |
| **Fmp bootstrapped** | 1.12% | 1.24%* | 1.24% | 0.97% | 1.48%* | 1.02% | 1.20% | 1.43% |
| **Modified q-sample V2** | 0.94% | 0.96% | 1.17% | 1.25% | 1.44%* | 1.88%* | 0.87% | 1.52% |
| **Modified q-sample V4** | 0.86% | 1.05% | 1.15% | 0.89% | 1.14% | 0.97% | 1.13% | 1.43% |

898
899

Chesnaye et al: Objective measures for detecting the ABR: comparisons in specificity, sensitivity, and detection time