1   **Title:** Objective detection of auditory steady-state responses based on mutual information: Receiver

2   operating characteristics and validation across modulation rates and levels

3

4   **Running head:** ASSR detection via information theory

5

6   **Authors:** Gavin M. **Bidelman**[1,2,3] and Claire **McElwain**[2]

7

8

9   **Affiliations:**

10  [1]Institute for Intelligent Systems, University of Memphis, Memphis, TN, USA, 38152

11  [2]School of Communication Sciences & Disorders, University of Memphis, Memphis, TN, 38152

12  [3]Univeristy of Tennessee Health Sciences Center, Department of Anatomy and Neurobiology, Memphis,

13  TN, 38163

14

15

16

17

18

19

20

21

22

23

24

25

26

27
28  **Address for editorial correspondence:**

29  Gavin M. Bidelman, PhD

30  School of Communication Sciences & Disorders

31  University of Memphis

32  4055 North Park Loop

33  Memphis, TN, 38152

34  TEL: (901) 678-5826

35  FAX: (901) 525-1282

36  EMAIL: g.bidelman@memphis.edu

1

37     **ABSTRACT**

38     Auditory steady-state responses (ASSRs) are sustained potentials used to assess the physiological

39     integrity of the auditory pathway and objectively estimate hearing thresholds. ASSRs are typically

40     analyzed using statistical procedures in order to remove the subjective bias of human operators. Knowing

41     when to terminate signal averaging in ASSR testing is also critical for making efficient clinical decisions

42     and obtaining high-quality data in empirical research.  Here, we investigated a new detection metric for

43     ASSRs based on mutual information (MI) [Bidelman, G. M. (2014). Objective information-theoretic

44     algorithm for detecting brainstem evoked responses to complex stimuli. *J. Am. Acad. Audiol.,* 25(8), 711-

45     722], previously bench tested using only a single suprathreshold stimulus. ASSRs were measured in n=10

46     normal hearing listeners to various stimuli varying in modulation rate (40, 80 Hz) and level (80 – 20 dB

47     SPL). MI-based classifiers applied to ASSRs recordings showed that accuracy of ASSR detection ranged

48     from ~75 - 99% and was better for 40 compared to 80 Hz responses and for higher compared to lower

49     stimulus levels. Detailed receiver operating characteristics (ROC) were used to establish normative ranges

50     for MI for reliable ASSR detection across levels and rates (MI=0.9-1.6). Relative to current statistics for

51     ASSR identification (F-test), MI was found to be a more efficient metric for determining the stopping

52     criterion for signal averaging. Our new results confirm that MI can be applied across a broad range of

53     ASSR stimuli and might offer improvements to conventional objective techniques for ASSR detection.

54

55

56     **Keywords:** Auditory evoked potentials (AEPs); auditory stead state response (ASSR); evoked potential

57     classification; *F*-test; objective audiometry

58

59

## INTRODUCTION

60

61      Auditory steady-state responses (ASSRs) are sustained evoked potentials typically elicited by

62      amplitude or frequency modulated signals. ASSRs offer a rapid physiological assessment of hearing

63      function and can be used to estimate full audiogram thresholds simultaneously in both ears (Cone-Wesson

64      et al., 2002; John & Picton, 2000; Picton et al., 1998). ASSRs are also preferred over other

65      electrophysiological measures (e.g., auditory brainstem response, ABR) because response detection is

66      based on a statistical comparison between signal and noise power in the evoked potential average rather

67      than human waveform inspection (Dobie & Wilson, 1996; John & Picton, 2000; Sturzebecher & Cebulla,

68      2013; Vidler & Parker, 2004). This objectivity is beneficial as it avoids subjective operator interpretation

69      and bias in determining the presence/absence of a response and quality of the auditory evoked potential

70      (AEP) recording (Bidelman, 2014; Bogaerts et al., 2009; Vidler & Parker, 2004).

71      Current approaches to analyze ASSRs typically involve frequency-domain measures where a

72      statistic is applied to the response spectrum in order to determine the significance of the signal's amplitude

73      relative to the surrounding noise floor (Dobie & Wilson, 1996; John & Picton, 2000; Sturzebecher &

74      Cebulla, 2013; Vidler & Parker, 2004). Several statistics have been proposed in the literature including the

75      $F$-test and magnitude-squared coherence (MSC) (Champlin, 1992; Dobie & Wilson, 1996). In all cases,

76      these statistics become more powerful with increasing number of trials. As such, a stopping rule can be

77      applied when a criterion value or significance level is achieved (e.g., $p<0.05$). Such metrics are currently

78      available in several commercial AEP systems. However, it remains unclear if these are the most optimal

79      statistics for characterizing sustained AEPs. Arguably, metrics like the $F$-test are somewhat limited

80      because they are usable only on specific features of the stimulus (e.g., power at the modulation frequency).

81      Consequently, these metrics cannot be broadly applied to sustained AEPs elicited by more complex sounds

82      (e.g., multi-frequency, time-varying stimuli) that have proven more useful in characterizing central

83      disorders of the auditory nervous system (e.g., Bidelman et al., 2017; Cone-Wesson et al., 2002; Johnson

84      et al., 2005; Purcell et al., 2004; Rocha-Muniz et al., 2012). Novel statistical approaches might offer higher

85      sensitivity and/or flexibility for detecting ASSRs and other sustained AEPs.

86      Towards this end, we have recently developed a new statistical method for detecting sustained

87      auditory potentials based on mutual information (MI) (Bidelman, 2014), a metric adopted from

88      information theory and image processing (for review, see Pluim et al., 2003). The essence of our approach

89      is to compare the spectrographic representations of the stimulus signal to that of the neural response

90      (Bidelman, 2014). MI enables us to characterize signal similarity by considering both linear and nonlinear

91      dependencies between neural responses and the evoking acoustic stimulus. By applying this metric to

92      signal and response spectrogram images, we take advantage of the full three-dimensional nature of the

93      AEP's time-frequency-amplitude information. In our previous bench tests, we showed that MI could

3

94　reliably detect speech-evoked frequency-following responses (FFRs) (Bidelman, 2014) and 40 Hz ASSRs

95　(Bidelman & Bhagat, 2016) from sham (EEG noise) recordings with ~90% accuracy. Moreover, we

96　reported that MI was superior to human observer judgements (Bidelman, 2014), was more robust in some

97　cases than the *MSC* and *F*-test (Bidelman & Bhagat, 2016), and yielded higher efficiency in detecting

98　ASSRs in shorter recording times than conventional statistical algorithms (Bidelman & Bhagat, 2016).

99　While promising, our previous investigations bench testing MI used only a *single* suprathreshold stimulus.

100　Thus, it remains unclear if MI can be more broadly applied to detect ASSRs elicited under a range of

101　stimulus parameters including different modulation rates and levels. Furthermore, the criterion threshold

102　for MI we used previously was estimated via computational modeling (Bidelman, 2014). Thus, it is not

103　clear from our previous studies whether this is the most appropriate criterion for detecting ASSRs evoked

104　by different stimulus levels and rates or if it was idiosyncratic to the one stimulus in our prior report.

105　Normative data reported here allowed us to address these open questions and recommend ranges for the

106　MI metric based on its performance (e.g., sensitivity) detecting a wider variety of ASSR responses.

107　Understanding the performance of MI detection across different stimulus settings is critical if the response

108　is to be eventually used for objective audiometry (Picton et al., 1998; John & Picton, 2000; Cone-Wesson

109　et al., 2002).

110　　　　The present study aimed to more fully characterize the performance of an MI-based classifier for

111　detecting ASSRs across a broader range of stimulus parameters. We assessed ASSR detection for

112　responses recorded at different modulation frequencies (40 Hz, 80 Hz) to assess the metric's dependence

113　on stimulus modulation rate (and thus putative site of the ASSR generator) (e.g., cortex vs. brainstem:

114　Herdman et al., 2002). Additionally, we parametrically varied stimulus level across a large dynamic range

115　(80–20 dB SPL) to evaluate the level-dependence of MI in detecting ASSRs. This latter manipulation is

116　important given the application of ASSRs for threshold estimation (Johnson & Brown, 2005;

117　Sturzebecher & Cebulla, 2013). Receiver operating characteristics (ROC) allowed us to characterize how

118　different choices of MI criterion values affect ASSR detection and thus, allowed us to established a

119　normative operating range for the metric and guide its future implementation. We further evaluated the

120　efficacy of the MI algorithm by comparing its application as a stopping criterion for signal ongoing

121　averaging against other "gold-standard" statistical approaches (i.e., F-test; John & Picton, 2000).

122　**METHODS & MATERIALS**

123　*Participants*

124　　Ten young, normal-hearing listeners (5 male, 5 female; age: 23.7±1.94 years) participated in the

125　experiment. All participants had normal hearing thresholds (≤ 15 dB HL, octave frequencies 250–8000

126　Hz) bilaterally, were right handed (Oldfield, 1971), and were native speakers of American English.

4

127 Participants gave written-informed consent in compliance with a protocol approved by the University of

128 Memphis Institutional Review Board (Protocol #2370).

129 *Stimuli*

130 ASSRs were evoked by sinusoidal amplitude modulated (SAM) tones with a carrier frequency ($f_c$) of

131 1000 Hz and modulation frequencies ($f_m$) of 40 Hz or 80 Hz (100% modulation depth). Stimulus duration

132 was 200 ms (including 5 ms onset/offset ramping to minimize onset components) following our previous

133 report (Bidelman & Bhagat, 2016). Stimuli were delivered binaurally via ER-2 insert earphones

134 (Etymotic Research) at levels of 80, 60, 40, and 20 dB SPL using alternating polarity. In addition to these

135 stimulus conditions, sham recordings were obtained by presenting stimuli with the inserts removed from

136 participants' ears (e.g., Aiken & Picton, 2008; Bidelman, 2014). Shams provided baseline, control

137 recordings of "neural noise" (Bidelman, 2014; Bidelman & Bhagat, 2016).

138 *Electrophysiological recordings*

139 ASSR recording procedures and stimuli were similar to our previous report (e.g., Bidelman &

140 Bhagat, 2016). EEGs were recorded between Ag/AgCl disc electrodes placed on the scalp at the high

141 forehead at the hairline, referenced to linked mastoids (A1/A2) (mid-forehead= ground). Interelectrode

142 impedances were ≤ 5 kΩ. Continuous EEG signals were digitized at 10 kHz (SynAmps RT amplifiers;

143 Compumedics Neuroscan). EEGs were windowed [0-200 ms], filtered (30-1000 Hz), and averaged in the

144 time domain to obtain ASSR waveforms for each stimulus. Listeners heard 2500 repetitions of the

145 stimulus token presented at an interstimulus interval of 5 ms. Post-processing and analyses were

146 performed using custom routines coded in MATLAB® 2015b. (The MathWorks, Inc.)

147 *Mutual information (MI) detection metric*

148 We computed the mutual information (MI) between spectrographic representations of the

149 stimulus and neural response to index the degree to which neural responses captured spectrotemporal

150 details of the acoustic input. Details of this metric are fully elaborated in our previous studies bench

151 testing this metric for AEP detection (Bidelman, 2014; Bidelman & Bhagat, 2016). MI is a dimensionless

152 quantity (measured in bits), which measures the degree of linear and nonlinear dependence between two

153 signals (*A* and *B*). In the specific case where *A* and *B* are two spectrograms, MI computes the *dependence*

154 or *similarity* between the two images (Pluim et al., 2003).

155 MI was computed between the stimulus and each neural response spectrogram allowing us to

156 assess the degree to which neural responses reflected spectrotemporal properties of the evoking stimulus

157 (Bidelman, 2014). Spectrograms were computed using the "spectrogram" routine in MATLAB and

158 converted to grayscale images. This routine computed a $2^{14}$ point FFT in consecutive 50 ms segments

5

159   (Hamming windowed) computed every 3 ms (Bidelman, 2014)[1]. Time waveforms were zero-padded to

160   minimize edge effects and ensure that spectrograms ran to the end of the signal's duration. Identical

161   parameters were used to compute both the stimulus and response spectrograms. SAM tone stimulus

162   spectrograms were squared prior to computing MI to account for the half-wave rectification that is

163   applied during the cochlear transduction process (Bidelman & Bhagat, 2016; Lins et al., 1995; Oxenham

164   et al., 2004).

165   *Receiver operating characteristics (ROC) for the MI classifier*

166        After determining a criterion (i.e., decision rule) for MI empirically from our data (see *Results*),

167   we then applied this threshold ($MI_\theta$) as a binary classifier to ASSR and sham recordings. Recordings

168   yielding $MI \geq MI_\theta$ were classified as neural responses whereas recordings with $MI < MI_\theta$ were considered

169   to be noise (i.e., no response) (Bidelman, 2014). Classifier performance was evaluated by computing

170   standard signal detection theory and ROC metrics including true and false positive rates. ROC analyses

171   also allowed us to validate the acceptable range of MI values that yielded above chance detection of

172   ASSRs from noise. For a given value of MI, sensitivity was computed as the percentage of actual ASSR

173   recordings correctly identified; false-positive rate as the percentage of sham recordings (i.e, "neural

174   noise") erroneously classified as a biological ASSR response. ROC curves were constructed for each

175   modulation rate (40 Hz, 80 Hz) and level (80 – 20 dB SPL) to characterize the overall performance of the

176   MI classifier across the different stimulus settings.

177   *Comparison of MI to the F-test*

178        To test the efficiency of MI as a stopping criterion for signal averaging, we computed MI on a

179   sweep-by-sweep basis as accumulating trials were added to the ongoing ASSR average. This was

180   repeated separately for each level and modulation rate. Similarly, we compared the "online" development

181   of MI against the well-known F-test (Dobie & Wilson, 1996; John & Picton, 2000) used in commercial

182   ASSR recording systems (e.g., Bio-logic MASTER II; Intelligent Hearing Systems SmartEP-ASSR).

183   While other detection metrics are available (e.g., MSC) we have previously shown that MI is most

184   comparable in detection performance to the F-test (MSC performs more poorly) (Bidelman & Bhagat,

185   2016), and thus, represents a stringent comparison to benchmark against. The underlying assumption of

186   this approach is that in the spectral domain, ASSR energy should be localized to a frequency bin near the

187   stimulus modulation frequency; activity in adjacent bins contain only random noise with zero mean and

---

[1] Window length changes the spectral resolution of the resulting spectrogram which could impact the computation of MI when comparing the stimulus and responses spectrograms. In initial analyses, we varied the sliding window length parametrically from 25 ms to 100 ms. However, in pilot testing, we found no appreciable changes in the accuracy of response detection for different window lengths (data not shown). Consequently, we adopted a 50 ms window, equivalent to a spectral resolution of 20 Hz. This is able to resolve both 40 Hz and 80 Hz components and is consistent with our previous studies (Bidelman, 2014; Bidelman & Bhagat, 2016).

188 variance distributed equally across the noise bins (John & Picton, 2000). The ratio of signal power to the

189 sum of the powers in $N$ adjacent frequency bins is distributed according to an $F$ distribution with 2 and

190 $2N$-1 degrees of freedom (John & Picton, 2000). In the current study, we used N=12 frequency bins

191 surrounding the target signal. We then compared our measured F-ratio against the critical F-value with 2

192 and 23 degrees of freedom and obtained a corresponding $p$-value for response detection. Traces yielding

193 $p<0.05$ were deemed to have response energy at the *fm* frequency that was significantly above the

194 surrounding noise floor. Comparison between the MI and F-test statistical metrics allowed us to relate

195 their performance and determine differences in their stopping rule for signal averaging, that is, the

196 number of trials where each measure detected the presence of ASSRs.


**RESULTS**

198 *ASSR responses*

199 ASSR time waveforms and spectra are shown for actual and sham recordings in Figure 1. Spectra

200 illustrate response energy at the modulation rates (40 Hz or 80 Hz) and their upper harmonics for ASSR

201 but not sham recordings (gray trace). These findings confirm that ASSRs contained robust phase-locked

202 neural activity whereas sham recordings contain no ASSR response (nor stimulus artifact) and are thus

203 suitable for use as "catch trials" in validating our MI detection metric (Bidelman & Bhagat, 2016). As

204 expected, ASSR amplitudes also decreased with decreasing stimulus level and were only weakly above

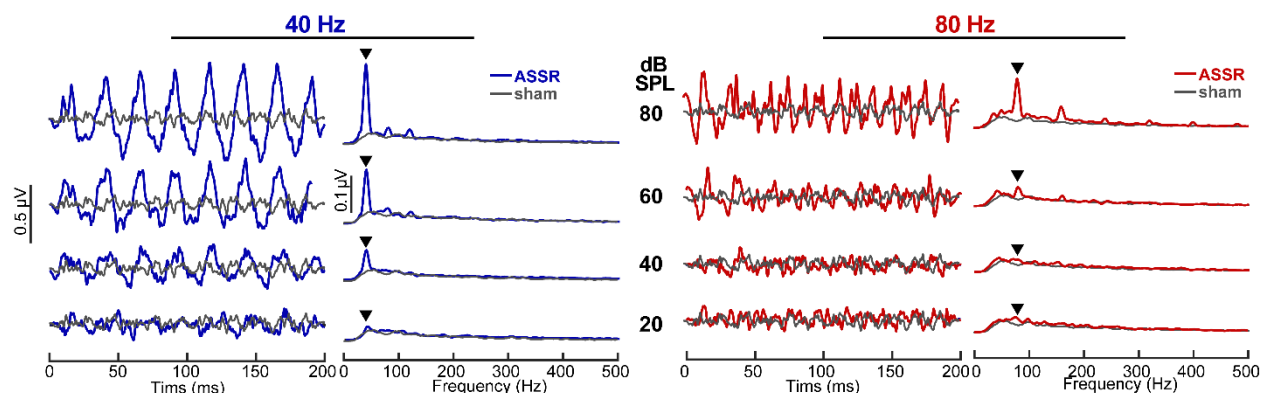205 the noise floor at 20 dB SPL.



**Figure 1:** Auditory steady-state response (ASSR) waveforms and spectra elicited by 40 Hz (*left*) and 80 Hz (*right*) SAM tones (*fc* = 1 kHz) showing the level dependence of responses. ASSR waveforms show phase locking at the stimulus modulation rate and first few harmonics which progressively weakens with decreasing level, approaching the noise floor at ~20 dB SPL. Gray traces, sham recording in which the earphone was removed from the ear canal (i.e., EEG noise floor). ▼=response energy at the *fm*.

206

207

208

7

209    *Performance and ROC characteristics of the MI classifier*

210         Examples of MI computed between the 40 Hz SAM stimulus and responses are shown for

211    different stimulus levels in Figure 2A. MI decreases at lower stimulus levels indicating weaker

212    dependence between the stimulus and ASSR response. At high intensities (80 dB SPL) ASSR

213    spectrograms show strong dependence on the evoking stimulus spectrogram and MI is large. Nearer

214    threshold (20 dB SPL), ASSRs are dominated by background EEG noise, implying that the averaged

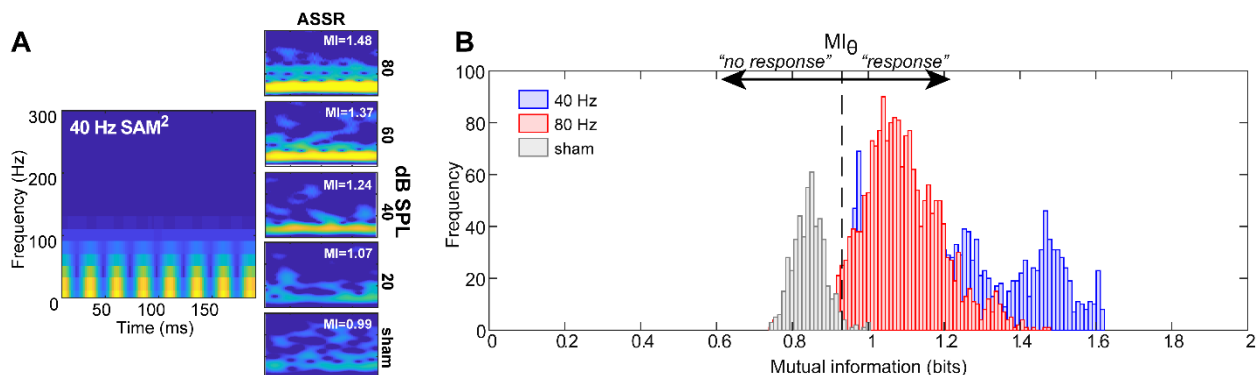215    neural response shares less information with the stimulus, which consequently yields a low MI[2].



**Figure 2:** Characterizing quality of ASSR recordings using mutual information (MI). (**A**) (*left*) Rectified stimulus spectrogram for a 40 Hz SAM tone stimulus. (*right*) Spectrograms of ASSRs recorded for a descending level series. The inset of each panel indicate the MI computed between each response spectrogram and that of the stimulus (Bidelman, 2014; Bidelman & Bhagat, 2016). With decreasing level, time-frequency representations of the neural ASSRs show less correspondence with that of the stimulus as indicated by decreasing values of MI. (**B**) Signal detection theory analysis to determine an optimal criterion (MI$_\theta$) for the MI response classifier. Shown here are the distribution (probability density functions) of MI values for 40 and 80 Hz ASSRs (pooled across stimulus levels) and sham recordings. MI is always larger for true vs. sham recordings. The criterion MI$_\theta$ = 0.93 segregates 95% of suprathreshold ASSRs from sham noise. From a classifier perspective, recordings containing an MI > MI$_\theta$ are predicted to contain a true neural ASSR response whereas recordings with MI < MI$_\theta$ are considered noise (no response).

216         Our first aim was to empirically determine a decision rule for MI for use in detecting ASSR

217    responses. To this end, signal detection theory was used to determine an optimal criterion (MI$_\theta$) for the

218    MI classifier from the recordings. Figure 2B shows the probability density functions of MI values for all

219    trials and subjects for the 40 Hz and 80 Hz ASSRs (pooling across levels) and sham recordings. On

220    average, MI values range from ~1 to 1.5 across all stimulus combinations. All ASSRs are, to varying

221    degrees, linearly separable along the MI decision axis compared to sham recordings which elicit weak MI

222    (~0.9). In the current study, MI$_\theta$ = 0.93 was taken as the criterion value because 95% of the data (i.e., MIs

223    for ASSR responses) fell above this threshold; consequently, the false positive rate was 5%. From a signal

224    detection standpoint, this implies that any arbitrary recording for which MI > MI$_\theta$ will predicted to

225    contain a true ASSR response whereas recordings with MI < MI$_\theta$ are considered noise (no response). MI$_\theta$

---

[2] Non-zero MI is observed even for sham recordings suggesting some shared-time-frequency information between the stimulus and neural noise. We attribute this to myogenic noise of the EEG which is strong for frequencies < 40 Hz. The SAM tone stimulus also has significant low frequency energy below < 40 Hz. Thus, even in the absence of a stimulus, spectral energy below < 40 Hz in both the stimulus and "neural noise" can produce a non-zero MI. This can be taken as the floor of the metric.

226     =0.93 was determined to be the optimal decision rule for ASSR detection and was used in subsequent

227     analyses.

228          Classifier performance of the MI metric is show in Figure 3 as ROC curves. Each panel

229     represents the true (sensitivity) vs. false positive (1- specificity) rate for distinguishing ASSRs from sham

230     recordings at different stimulus levels. The bowing of the ROC curve toward the upper left corner is

231     indicative of robust sensitivity in segregating signal from noise (i.e., higher $d$-prime). Each individual data

232     point represents the true/false positive rate for a different choice of $MI_\theta$. A criterion located at the

233     maximum curvature of the ROC curve represents the optimal decision rule for classification, one which

234     produces the highest sensitivity while minimizing false-positive detection (i.e., erroneously labeling a
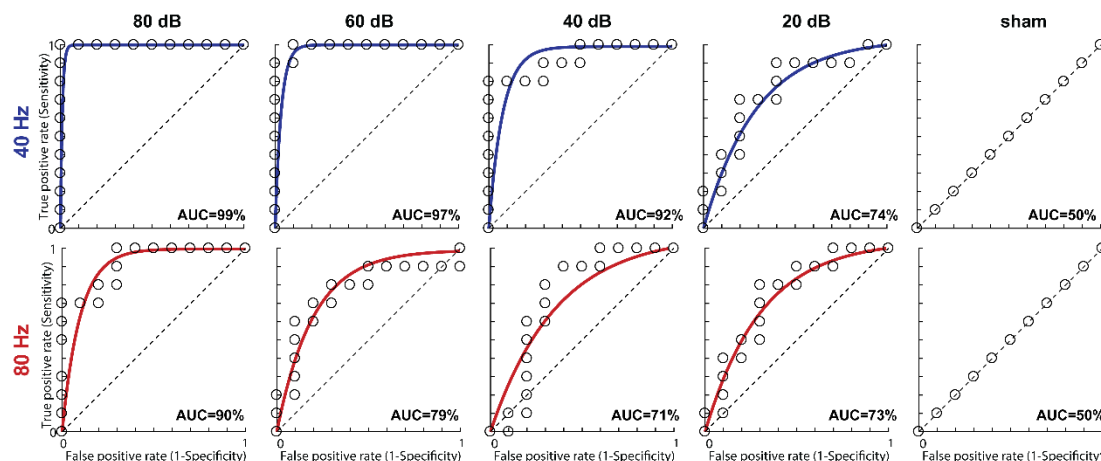
235     noise recording as an ASSR).



**Figure 3:** Receiver operating characteristic (ROC) curves for ASSRs evoked by different modulation rates and levels. Individual points denote the true positive (sensitivity) vs. false positive (1- specificity) rates for various values of MI in distinguishing true from sham recordings based on 2500 sweeps. Optimal sensitivity/specificity for the empirically derived criterion value ($MI_\theta$ =0.93) is repented by the point in the upper left corners of each ROC. Dotted lines correspond with chance performance (i.e., $d' = 0$). With $MI_\theta$ =0. 93, classification accuracy (AUC) is 99% for 40 Hz ASSRs at 80 dB SPL and 90% for 80 Hz ASSRs. Classification accuracy is better for low (40 Hz) compared to high (80 Hz) modulation rates and high vs. low level stimuli.

236

237          With decreasing levels, ASSRs become more difficult to segregate from EEG noise, as evident by

238     the ROC curves approaching chance performance (dotted lines) at 20 dB SPL. Overall classification

239     accuracy for the 40 Hz responses is near ceiling (99%) at 80 dB SPL, as indicated by the area under the

240     curve (AUC) (Hanley & McNeil, 1983). Classification accuracy weakens with decreasing level indicating

241     discriminating ASSRs from noise is more difficult nearer to threshold. Nevertheless, classification

242     remains high (74%) for the 40 Hz responses at 20 dB SPL. Accuracy in detecting 80 Hz responses is 10-

243     15% poorer compared to 40 Hz responses but still remains well above chance (73%) even at the lowest

244     intensity tested. These operating characteristics demonstrate that the MI between a stimulus and neural

245     response provides an objective means for detecting ASSRs across various levels and modulation rates.

9

246    *Acceptable ranges of MI for ASSR detection*

247        While the empirically derived criterion $MI_\theta = 0.93$ represents the *optimal* threshold for detecting

248    ASSRs (5% false positive), our ROC characterizations reveal there is a *range* of acceptable MI values

249    that could be used to reliably detect neural responses. Figure 4 shows the overall accuracy of detecting

250    ASSRs from shams for different choices of MI for 40 Hz (Fig. 4A) and 80 Hz (Fig. 4B) responses. Each

251    family of functions shows the overall accuracy in correctly detecting ASSRs from noise using different

252    MI cutoffs. The reduction in peak accuracy across curves indicates a level-dependent effect in

253    classification accuracy. Consistent with ROC results, MI is less robust at detecting ASSRs evoked by

254    weaker stimulus levels. Nevertheless, there is a *range* of MI values that still allow above-chance detection

255    of the response (Fig. 4C). MI ranges were extracted from the width of each level-dependent accuracy

256    function shown in panels A and B and show the acceptable range of MI cutoff thresholds that allow from

257    above-chance detection. For the 40 Hz response, acceptable values of MI range from 0.9–1.6 for high

258    level (80 dB) stimuli. This allowable range is reduced with decreasing level; 40 Hz responses are

259    detectable at 20 dB with MI values between 0.9–1.3. Similar results were obtained for the 80 Hz

260    responses, although the acceptable MI range was reduced at both high- (0.9-1.4) and low-level (0.9-1.2)
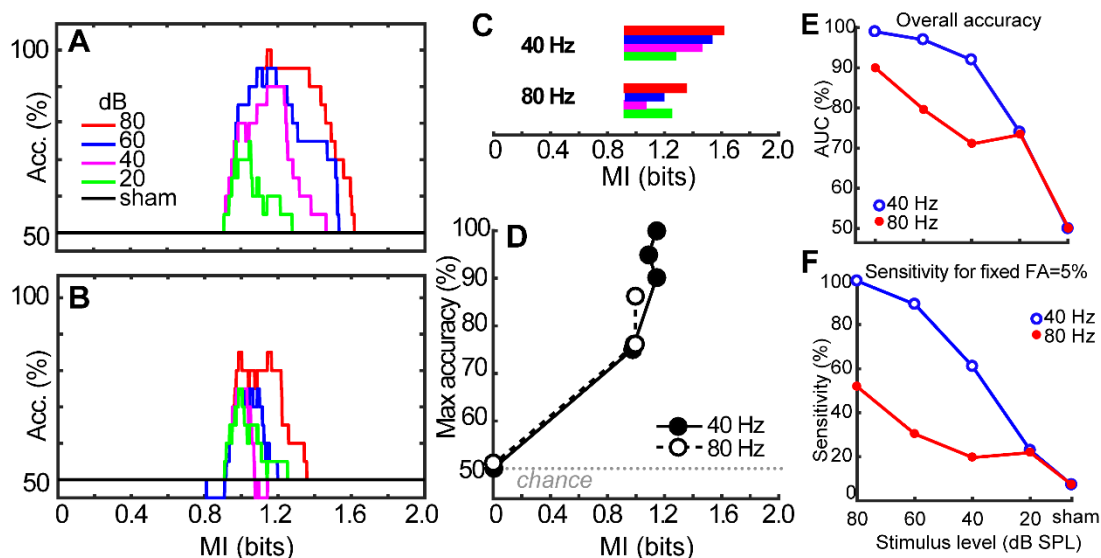
261    stimuli.



**Figure 4:** Acceptable MI values for detecting suprathreshold ASSRs across stimulus level and modulation rates. (**A-B**) Level-dependent classification accuracy for the 40 Hz (A) and 80 Hz (B) responses. Each family of functions shows the overall accuracy in correctly detecting ASSRs from noise using difference MI cutoffs. (**C**) Range of acceptable MI values for classifying 40 and 80 Hz ASSRs above chance levels. Ranges were extracted from the width of each level-dependent accuracy function shown in panels A and B. (**D**) Max accuracy for detecting ASSRs and the corresponding MI. Max accuracy was extracted from the peak of each level-dependent accuracy function of panel A and B. Note that some points for the 80 Hz responses overlap. (**E**) Overall accuracy across stimulus levels and modulation rates. Accuracies were extracted from ROC functions (e.g., Fig. 4) as the area under the curve (AUC). (**F**) Sensitivity of the MI metric controlling (fixing) false positive rate at 5%. Accuracy and sensitivity are better for 40 Hz compared to 80 Hz responses, decrease with decreasing stimulus level, but remain well above chance.

262

10

263   MI values corresponding to maximum classification accuracy (i.e., peak of functions in Figs. 4A-

264 B) are shown in Fig. 4D. Maximum accuracy is obtained with an MI≈1 (cf. $MI_\theta$). Collectively, these

265 results help provide a normative tolerance range and optimal choice of MI values for using it as an ASSR

266 classifier.

267   Typically, the performance of a diagnostic or detection method is evaluated by considering the

268 sensitivity and specificity of the measure. However, it is also useful to evaluate a diagnostic's true

269 positive rate (i.e., sensitivity) for a *fixed* false positive rate (e.g., 5%). Figure 4E-F shows the overall

270 accuracy of the MI metric (AUC) and sensitivity at a fixed 5% false positive rate (i.e., 95% specificity)

271 across stimulus levels and modulation rates. For 40 Hz ASSRs, performance ranges from 100/95%

272 sensitivity/specificity at 80 dB SPL to 20/95% sensitivity/specificity at 20 dB SPL. For 80 Hz ASSRs,

273 performance ranges from 55/95% sensitivity/specificity at 80 dB SPL to 20/95% sensitivity/specificity at

274 20 dB SPL.

275 *MI as a criterion for terminating signal averaging*

276   In addition to detection, an objective metric should be suitable as a stopping criterion for online

277 signal averaging. Figure 5 shows the growth in MI (present study; Bidelman, 2014; Bidelman & Bhagat,

278 2016) and conventional F-test (Dobie & Wilson, 1996; John & Picton, 2000) metric during ASSR

279 recordings as a function of the number of trials in the ongoing average. In general, each metric improves

280 with additional trials and asymptotes as the running AEP stabilizes. Response growth is faster for 40 Hz

281 relative to 80 Hz responses and at higher (80 dB SPL) compared to lower (20 dB SPL) intensities. For

282 high-level 40 Hz ASSRs, responses exceed the MI and F-test stopping criteria (MI=0.9; F-test: p=0.05)

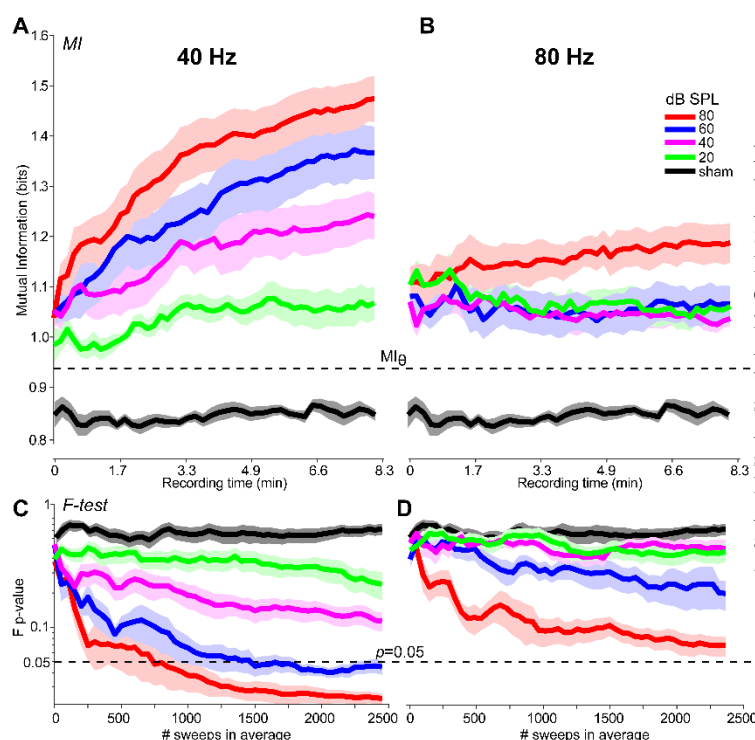283 by ~50 and 750 sweeps, respectively, corresponding to < 1min (MI) vs. 2.5 min (F-test) of recording time.



**Figure 5:** Comparison of the growth in MI to the *F*-test during online ASSR recording. Sweep-by-sweep ASSR detection based on MI for the 40 Hz (**A**) and 80 Hz (**B**) responses at each level. (**C-D**) Improvement in ASSR detection using the conventional F-test procedure (John & Picton, 2000). Dotted lines denote the criterion threshold for ASSR detection under each metric ($MI_\theta$: 0.93; F-test: *p*<0.05). Abscisse show both the number of sweeps and corresponding recording time for online ASSR averaging. 40 Hz responses exceed the MI stopping criteria within ~50 sweeps (< 1 min); longer recording times are needed for detecting ASSRs using the F-test [e.g., ~750 sweeps (2.5 min) are required for the 40 Hz reponse@80 dB SPL and ~1500 sweeps at 60 dB SPL]. Shading = ±1 s.e.m.

11

284   More extended recording durations (sweeps) are needed for detecting low-level ASSRs and the 80 Hz

285   responses, which sometimes do not achieve the criterion thresholds (e.g., Fig. 6D). As an expected

286   control, MI remains invariant sweep-to-sweep for sham (noise) recordings.

287   **DISCUSSION**

288         In the current study, we demonstrate the tolerance of a new, objective statistical approach to

289   detect ASSRs based on mutual information (Bidelman, 2014). The technique quantifies the quality of

290   ASSRs by considering the linear and nonlinear dependences between the rich time-frequency information

291   provided by the signal and response spectrograms. Our previous reports bench testing the MI classifier

292   demonstrated its superiority over "gold standard" judgments of human observers (Bidelman, 2014) and

293   other objective techniques for ASSR detection (e.g., MSC, *F*-test) (Bidelman & Bhagat, 2016) for

294   suprathreshold (70-80 dB SPL) stimuli.  Here, we extend these previous findings by showing that MI can

295   be used for response detection across a broader range of ASSR-evoking stimuli including different

296   combinations of levels and modulation rates.

297         Overall performance accuracy in distinguishing true neurobiological responses from noise using

298   our MI metric was >90% for high level stimuli (80 dB SPL) and remained well-above chance (73%) for

299   levels nearer to threshold (20 dB SPL). More importantly, our results establish a normative tolerance

300   range for MI criterion values (MI = 0.9 − 1.6) that allow for robust detection of ASSRs across different

301   modulation rates and intensities. However, as determined by ROC analyses, the most optimal

302   classification of ASSRs is achieved with a criterion $MI_\theta = 0.93$. Lastly, we showed that MI increases

303   monotonically with increasing number of stimulus presentations (i.e., trials) and can, for some stimulus

304   conditions, detect ASSRs in a fewer number of trials compared to conventional ASSR detection

305   procedures (i.e., F-test; Dobie & Wilson, 1996; John & Picton, 2000).

306         In prior studies, we previously showed that MI could be applied to other classes of AEPs

307   including speech-evoked FFRs (Bidelman, 2014) as well as high-level ASSRs (Bidelman & Bhagat,

308   2016). MI is an information-theoretic measure that is "distribution free" and therefore requires fewer

309   assumptions than other statistical approaches (e.g., *F*-test)*, which utilize parametric (distribution-based)

310   statistics. Unlike other metrics, MI can also be easily applied to time-varying signals (Bidelman, 2014).

311   Thus, in addition to potentially broader application, MI may offer a useful alternative to other ASSR

312   response detection approaches currently employed in commercial hardware.

313         The more comprehensive stimulus set in this compared to our previous studies (Bidelman, 2014;

314   Bidelman & Bhagat, 2016) allows for a more comprehensive characterization of MI's effectiveness as an

315   ASSR classifier. Several observations are worth noting regarding the metric's performance. First, while

316   MI can successfully detect the presence of ASSRs at different modulation rates (Fig. 3), we found overall

12

accuracy was generally higher for 40 Hz compared to 80 Hz responses. Thus, while MI can successfully

detect ASSRs across a wide range of stimulus levels and modulation rates, it is more accurate and

sensitive for 40 Hz responses and higher, compared to lower level stimuli. The more optimal performance

at 40 Hz is likely due to the higher signal-to-noise ratios and more robust amplitudes of ASSRs to low vs.

high-frequency modulation rates (present study, Fig. 2; Galambos et al., 1981; Korczak et al., 2012;

Purcell et al., 2004). Indeed, by early adolescence, the 40 Hz response is nearly twice the amplitude of the

80 Hz response (Pethe et al., 2004). Moreover, unlike their 80 Hz counterparts, 40 Hz responses are

highly dependent on subject state: low *fm* responses are reliably recorded only in awake individuals

(Cohen et al., 1991; Korczak et al., 2012; Kuwada et al., 1986) and are eradicated with anesthesia

(Galambos et al., 1981; Kuwada et al., 2002). These properties have limited the utility of the 40 Hz ASSR

for infant testing. Thus, while we have confirmed that MI is efficacious for detecting suprathreshold

ASSRs across different modulation rates, MI would be less appropriate to monitor response detection for

low-level stimuli. This may limit the metric's utility for hearing threshold testing. Nevertheless,

suprathreshold ASSRs and other sustained AEPs do find clinical use [e.g., newborn hearing screenings

(American Academy of Pediatrics, 2007)]. Research applications typically involve complex paradigms,

multiple subject cohorts, and longer testing protocols. Our results therefore suggest that MI could offer a

means to collect sustained ASSR/AEP data in a more time-optimized manner and reduce valuable

recording time (present study; Bidelman, 2014; Bidelman & Bhagat, 2016).

Secondly, we find that MI has a smaller useable range (Fig. 4C) and lower accuracy/sensitivity

(Fig. 4D) for low-level, 80 Hz stimuli. This would tend the limit the metric's application for threshold

testing (Picton et al., 2005), particularly in infants (Stroebel et al., 2007). Additionally, neural generators

of the ASSR are dependent on the frequency of the stimulus modulation rate; high frequencies (80 Hz)

evoke brainstem generators whereas low-frequencies (40 Hz) recruit cortical sources (Herdman et al.,

2002; Kuwada et al., 2002). Thus, the fact that we observe superior performance for 40 Hz stimuli across

the board implies that MI might be more useful for monitoring cortical rather than subcortical neural

activity.

Lastly, sweep-by-sweep tracking of MI confirmed the metrics' efficiency as a stopping rule for

ASSR signal averaging. In this regard, we found that MI was able to detect 40 Hz ASSRs within ~1 min,

corresponding to < 50 stimulus trials. In contrast, using the F-test required considerably more stimulus

presentations; ~750 sweeps (2.5 min of testing) were needed to detect the 40 Hz response at 80 dB SPL

and ~1500 sweeps at 60 dB SPL. Moreover, our 80 Hz ASSRs never achieved the F-test criterion,

indicating that more than 2500 trials would be needed to detect those responses. While our data show that

that MI can offer a more efficient stopping rule for terminating averaging compared to the *F*-test, from a

practical standpoint, this improvement in testing time (1-2 min) is probably negligible. Nevertheless, our

13

351    data indicate that under some stimulus conditions, MI can detect ASSRs in half the number of trials (i.e.,

352    twice as efficient) as the gold-standard $F$-test.

353         In conclusion, the application of the MI metric to electrical response audiometry may provide

354    clinicians and researchers with a more robust tool to objectively evaluate the presence and quality of

355    sustained auditory AEPs. Calculation of MI could be easily incorporated into most commercially

356    available AEP systems similar to other statistical detection metrics already in place (e.g., *F-test, $F_{sp}$,*

357    *MSC*). Future studies are warranted to assess the performance of MI in infant and threshold ASSR testing.

358

359

14

## References

360

361  Aiken, S.J., Picton, T.W. 2008. Envelope and spectral frequency-following responses to vowel sounds.
362         Hear. Res. 245, 35-47.

363  American Academy of Pediatrics, J.C.o.I.H. 2007. Position Statement: Principles and guidelines for early
364         hearing detection and identification. Pediatrics 120, 898-921.
365  Bidelman, G.M. 2014. Objective information-theoretic algorithm for detecting brainstem evoked
366         responses to complex stimuli. J. Am. Acad. Audiol. 25, 711-722.
367  Bidelman, G.M., Bhagat, S.P. 2016. Objective detection of auditory steady-state evoked potentials based
368         on mutual information. Int. J. Audiol. 55, 313-319.
369  Bidelman, G.M., Lowther, J.E., Tak, S.H., Alain, C. 2017. Mild cognitive impairment is characterized by
370         deficient hierarchical speech coding between auditory brainstem and cortex. J. Neurosci. 37,
371         3610-3620.
372  Bogaerts, S., Clements, J.D., Sullivan, J.M., Oleskevich, S. 2009. Automated threshold detection for
373         auditory brainstem responses: Comparison with visual estimation in a stem cell transplantation
374         study. BMC Neurosci. 10, 1-7.
375  Champlin, C.A. 1992. Methods for detecting auditory steady-state potentials recorded from humans.
376         Hear. Res. 58, 63-69.
377  Cohen, L.T., Rickards, F.W., Clark, G.M. 1991. A comparison of steady-state evoked potentials to
378         modulated tones in awake and sleeping humans. J. Acoust. Soc. Am. 90, 2467-79.
379  Cone-Wesson, B., Dowell, R.C., Tomlin, D., Rance, G., Ming, W.J. 2002. The auditory steady-state
380         response: comparisons with the auditory brainstem response. J. Am. Acad. Audiol. 13, 173-87.
381  Dobie, R.A., Wilson, M.J. 1996. A comparison of t test, F test, and coherence methods of detecting
382         steady-state auditory-evoked potentials, distortion-product otoacoustic emissions, or other
383         sinusoids. J. Acoust. Soc. Am. 100, 2236-2246.
384  Galambos, R., Makeig, S., Talmachoff, P. 1981. A 40-Hz auditory potential recorded from the human
385         scalp. Proc. Natl. Acad. Sci. USA 78, 2643-2647.
386  Hanley, J.A., McNeil, B.J. 1983. A method of comparing the areas under receiver operating characteristic
387         curves derived from the same cases. Radiology 148, 839-43.
388  Herdman, A.T., Lins, O., van Roon, P., Stapells, D.R., Scherg, M., Picton, T. 2002. Intracerebral sources
389         of human auditory steady-state responses. Brain Topogr. 15, 69-86.
390  John, M.S., Picton, T.W. 2000. MASTER: A Windows program for recording multiple auditory steady-
391         state responses. Comput. Methods Programs Biomed. 61, 125–150.
392  Johnson, K.L., Nicol, T.G., Kraus, N. 2005. Brain stem response to speech: A biological marker of
393         auditory processing. Ear Hear. 26, 424-34.
394  Johnson, T.A., Brown, C.J. 2005. Threshold prediction using the auditory steady-state response and the
395         tone burst auditory brain stem response: a within-subject comparison. Ear Hear. 26, 559-76.
396  Korczak, P., Smart, J., Delgado, R., Strobel, T.M., Bradford, C. 2012. Auditory steady-state responses. J.
397         Am. Acad. Audiol. 23, 146-70.
398  Kuwada, S., Batra, R., Maher, V.L. 1986. Scalp potentials of normal and hearing-impaired subjects in
399         response to sinusoidally amplitude-modulated tones. Hear. Res. 21, 179-92.
400  Kuwada, S., Anderson, J.S., Batra, R., Fitzpatrick, D.C., Teissier, N., D'Angelo, W.R. 2002. Sources of
401         the scalp-recorded amplitude-modulation following response. J. Am. Acad. Audiol. 13, 188-204.
402  Lins, O.G., Picton, P.E., Picton, T.W., Champagn, S.C., Durieux-Smith, A. 1995. Auditory steady-state
403         responses to tones amplitude-modulated at 80-110 Hz. J. Acoust. Soc. Am. 97, 3051–3063.
404  Oldfield, R.C. 1971. The assessment and analysis of handedness: The Edinburgh inventory.
405         Neuropsychologia 9, 97-113.
406  Oxenham, A.J., Bernstein, J.G.W., Penagos, H. 2004. Correct tonotopic representation is necessary for
407         complex pitch perception. Proc. Natl. Acad. Sci. USA 101, 1421-1425.

408 Pethe, J., Muhler, R., Siewert, K., von Specht, H. 2004. Near-threshold recordings of amplitude
409      modulation following responses (AMFR) in children of different ages. Int. J. Audiol. 43, 339-45.
410 Picton, T.W., Dimitrijevic, A., Perez-Abalo, M.C., Van Roon, P. 2005. Estimating audiometric thresholds
411      using auditory steady-state responses. J. Am. Acad. Audiol. 16, 140-56.
412 Picton, T.W., Durieux-Smith, A., Champagne, S.C., Whittingham, J., Moran, L.M., Giguere, C.,
413      Beauregard, Y. 1998. Objective evaluation of aided thresholds using auditory steady-state
414      responses. Journal of American Academy of Audiology 9, 315-331.
415 Pluim, J.P., Maintz, J.B., Viergever, M.A. 2003. Mutual-information-based registration of medical
416      images: A survey. IEEE Trans. Med. Imaging 22, 986-1004.
417 Purcell, D.W., John, S.M., Schneider, B.A., Picton, T.W. 2004. Human temporal auditory acuity as
418      assessed by envelope following responses. J. Acoust. Soc. Am. 116, 3581-93.
419 Rocha-Muniz, C.N., Befi-Lopes, D.M., Schochat, E. 2012. Investigation of auditory processing disorder
420      and language impairment using the speech-evoked auditory brainstem response. Hear. Res. 294,
421      143-152.
422 Stroebel, D., Swanepoel, W., Groenewald, E. 2007. Aided auditory steady-state responses in infants. Int.
423      J. Audiol. 46, 287-292.
424 Sturzebecher, E., Cebulla, M. 2013. Automated auditory response detection: Improvement of the
425      statistical test strategy. Int J Audiol 52, 861-4.
426 Vidler, M., Parker, D. 2004. Auditory brainstem response threshold estimation: Subjective threshold
427      estimation by experienced clinicians in a computer simulation of a clinical test. Int. J. Audiol. 43,
428      417-429.

429

430