

# Intrasubject variability of early markers of sensorineural hearing damage

Nele De Poortere<sup>1</sup>, Sarineh Keshishzadeh<sup>2</sup>, Hannah Keppler<sup>1,4</sup>, Ingeborg Dhooge<sup>3,4</sup>, Sarah Verhulst<sup>2</sup>

<sup>1</sup>Ghent University, Department of Rehabilitation Sciences – Audiology

<sup>2</sup>Ghent University, Dept. of Information Technology – Hearing Technology @ WAVES

<sup>3</sup>Ghent University Hospital, Department of Ear, Nose and Throat, Belgium

<sup>4</sup>Ghent University, Department of Head and Skin

## Financial disclosure

Ghent University holds a patent on the RAM-EFR stimulation protocol (U.S. patent App. 17/791,985 (Inventors: Sarah Verhulst, Viacheslav Vasilkov). This work was supported by UGent BOF-IOP project: Portable Hearing Diagnostics: Monitoring of Auditory-nerve Integrity after Noise Exposure (EarDiMon), European research council proof of concept grant CochSyn (899858) and EU Innovation council Grant EarDiTech (101058278).

## Conflict of interest

The authors report there are no competing interests to declare.

## Data Availability Statement

The datasets generated during and/or analyzed during the current study are not publicly available due to ethical restriction but are available from the corresponding author on reasonable request.

## Acknowledgement:

The authors would like to thank Eef De Wilde and Ellen Sabau for their significant contribution in the data-collection.

## All correspondence should be addressed to:

Nele De Poortere, Faculty of Medicine and Health Sciences, Department of Rehabilitation Sciences, Corneel Heymanslaan 10 (2P1), Ghent University, 9000 Ghent, Belgium; nele.depoortere@ugent.be; +32 479 02 32 62.

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

31 ABSTRACT

32 The pursuit of noninvasive early markers for sensorineural hearing loss (SNHL) has yielded diverse  
33 measures of interest. However, comprehensive studies evaluating the test-retest reliability of multiple  
34 measures and stimuli within a single study are scarce, and a standardized clinical protocol for robust  
35 SNHL-markers remains elusive. To address these gaps, this study covers the intra-subject variability of  
36 potential EEG-biomarkers for cochlear synaptopathy (CS) and other SNHL-markers to determine their  
37 clinical suitability. Fifteen normal-hearing young adults underwent repeated measures of (extended high-  
38 frequency) pure-tone audiometry, speech-in-noise intelligibility, distortion-product otoacoustic emissions  
39 (DPOAEs), and auditory evoked potentials; comprising envelope following responses (EFR) and auditory  
40 brainstem responses (ABR). Results show high reliability in pure-tone audiometry, whereas the matrix  
41 sentence-test showed a significant learning effect. DP-grams and input-output functions' reliability varied  
42 across three evaluation methods with distinct SNR-based criteria for DPOAE-datapoints. EFRs  
43 demonstrated superior test-retest reliability compared to ABR-amplitudes. Our findings underscore  
44 careful interpretation of presumed noninvasive SNHL measures. While we confirm the robustness of  
45 tonal-audiometry, we found a confounding learning effect in longitudinal speech audiometry. DPOAE  
46 variability underscores the need for consistent ear probe replacement and meticulous measurement  
47 techniques and renders I/O-functions unsuitable for clinical application. As potential EEG-biomarkers of  
48 CS, EFRs are favored over ABR-amplitudes.

49 **Keywords:** Intrasubject variability – noninvasive early markers of sensorineural hearing loss-  
50 cochlear synaptopathy – normal hearing young adults – envelope followings response

51  
52  
53  
54  
55  
56  
57  
58

## 59 INTRODUCTION

60 In the realm of clinical application, the standard procedure for assessing auditory thresholds  
61 relies on conventional pure-tone audiometry. Consequently, studies targeting the evaluation of  
62 hearing damage related to aging, ototoxicity and excessive noise exposure have primarily  
63 centered on identifying permanent hearing threshold changes within the frequency range of  
64 250 to 8000 Hz (Cruickshanks et al., 2010; Rabinowitz et al., 2006). Temporary threshold shifts  
65 (TTS) resulting from noise exposure were historically regarded as less concerning markers for  
66 permanent hearing damage, as indicated by the National Institute of Occupational Safety and  
67 Health (NIOSH) in 1998. However, recent rodent studies have challenged this perspective,  
68 revealing that a noise-exposure-induced TTS coincides with permanent deficits at the synaptic  
69 level, a phenomenon referred to as cochlear synaptopathy (CS). CS predominantly affects the  
70 connections between type-I auditory nerve fiber terminals and inner hair cells (IHCs) (Furman  
71 et al., 2013; Kujawa & Liberman, 2009) and can result in supra-threshold temporal coding  
72 deficits (Bharadwaj et al., 2015). Unfortunately, CS mostly affects supra-threshold sound  
73 coding without affecting routine clinical hearing sensitivity measures such as the tonal  
74 audiogram or distortion-product otoacoustic emissions (DPOAEs) (Furman et al., 2013;  
75 Lobarinas et al., 2013). Nonetheless, CS is believed to contribute to symptoms such as tinnitus  
76 and hyperacusis, and is thought to underlie other perceptual challenges, including difficulties  
77 in discriminating sounds in complex acoustic environments and impaired temporal processing  
78 of sound and speech intelligibility. (Mehraei et al., 2016; Mepani et al., 2021; Plack et al., 2014;  
79 Schaette & McAlpine, 2011; Verhulst et al., 2018).

80 Hence, researcher persists in their quest for robust and sensitive objective measures for routine  
81 clinical applications. Simultaneously, the translation of recent research methods into clinical  
82 practice, for improved SNHL diagnosis and monitoring, faces challenges, particularly in test-  
83 retest reliability. Lastly, as novel SNHL-treatments are being developed, the need for robust

84 biomarkers of SNHL are necessary to monitor treatment benefits. Consequently, this study  
85 endeavors to explore the intra-subject variability of potential EEG-biomarkers for CS and other  
86 early SNHL-indicators within the same individuals, with the goal of assessing their suitability  
87 for clinical use.

### 88 **Early Indicators of Outer Hair Cell Loss**

89 According to literature, extended high frequency (EHF) audiometry, presents a more sensitive  
90 approach for the early detection of noise induced hearing loss (NIHL) compared to  
91 conventional frequencies (Wang et al. (2000), Lopes et al. (2009), and Singh et al. (2009).  
92 Additionally, EHF audiometry has proven to be a valuable predictive tool for identifying the  
93 risk of NIHL (Mehrparvar et al., 2011), highlighting its clinical significance in the field of  
94 hearing health assessment.

95 An alternative, or complementary, approach to behavioral audiometry involves DPOAEs,  
96 which closely reflect the integrity of cochlear structures; particularly the outer hair cells (OHC)  
97 (Jansen et al., 2009). Furthermore, DPOAEs are recognized for their sensitivity in detecting  
98 subtle cochlear damage before it manifests in pure-tone hearing threshold elevations (Coradini  
99 et al., 2007; Knight et al., 2007; Reavis et al., 2015). Moreover, DPOAEs offer distinct  
100 advantages such as rapid acquisition, non-participatory nature, and suitability for measurement  
101 by non-specialist personnel trained in emission assessments (Reavis et al., 2015). However,  
102 despite the crucial role of DPOAEs in audiology, there remains a noticeable scarcity of studies  
103 that concurrently investigate the most suitable DPOAE evaluation methodologies for clinical  
104 applications, along with the reliability of auditory thresholds, DP-amplitudes, and DP-  
105 thresholds within the same cohort of subjects.

106

107

## 108 **Indicators of Cochlear Synaptopathy**

109 Among the promising EEG biomarkers for CS are gross suprathreshold neural potentials, such  
110 as the peak I amplitude of the auditory brainstem response (ABR) and the strength of the  
111 Envelope Following Response (EFR) (Guest et al., 2017; Liberman et al., 2016). Animal  
112 models have highlighted the significance of the ABR-amplitude as a clinical metric for  
113 diagnosing different types of hearing damage, including noise-induced CS (Kujawa &  
114 Liberman, 2009; Skoe & Tufts, 2018). However, translating the utility of ABR amplitudes as  
115 early biomarkers from rodents to humans poses challenges. Firstly, intersubject variabilities  
116 arising from differences in head size and sex (Mitchell et al., 1989), individual variation in  
117 cochlear dispersion (Don et al., 1994), electrode resistance, and various sources of electrical  
118 noise between individuals or sessions (Plack et al., 2016), act as confounding factors that hinder  
119 its diagnostic utility in humans. Secondly, while the early wave I of the ABR is assumed to  
120 have diagnostic potential in individual listeners, its evaluation in relation to CS is still an  
121 exploratory area in human research studies. The sensitivity of ABR measurements to low  
122 spontaneous rate (SR) auditory nerve fibers, which are particularly vulnerable to CS, has been  
123 questioned due to the delayed onset response of these fibers (Bourien et al., 2014). ABRs are  
124 namely evoked by transient stimuli and primarily reflect onset responses, which tend to be  
125 relatively small in low SR fibers (Rhode & Smith, 1985; Taberner & Liberman, 2005). In this  
126 respect, EFR-strengths, defined as the summation of the signal-to-noise spectral magnitude at  
127 the fundamental frequency of the stimulus and its following three harmonics (Vasilkov et al.,  
128 2021), are promising. Unlike ABRs, the synchronized firing of ANFs in response to EFR  
129 stimuli with shallow modulation depths are believed to be primarily driven by low SR-ANFs,  
130 which are more sensitive to CS as they exhibit greater synchronization compared to high SR-  
131 fibers when responding to sinusoidal amplitude modulated (SAM) tones (Bharadwaj et al.,  
132 2014).

133 Given that CS has been assumed to impact suprathreshold hearing sensitivity and speech  
134 recognition abilities, particularly in challenging listening conditions (Kujawa & Liberman,  
135 2009; Lin et al., 2011; Parthasarathy & Kujawa, 2018; Skoe et al., 2019), various speech  
136 recognition in noise tests have been employed in human studies to explore CS (Grinn et al.,  
137 2017; Guest et al., 2018; Maele et al., 2021). Nonetheless, research into the within-subject  
138 variability of these tests has remained limited.

139 In sum, the reliability of proposed noninvasive metrics as early markers of SNHL in humans  
140 remains uncertain due to the challenges posed by both intra- and inter-subject variabilities,  
141 hindering their uptake within clinical practice. This study investigates the intra-subject  
142 variability of potential EEG biomarkers for CS and other early SNHL indicators within the  
143 same individuals across three sessions in a cohort of normal hearing listeners, to suggest which  
144 measures are most reliable for use in clinical diagnostics and SNHL monitoring.

## 145 MATERIALS AND METHODS

### 146 PARTICIPANTS AND STUDY DESIGN

147 Fifteen young adults, nine men and six women, aged between 18 and 25 years (mean age 21.0  
148 years  $\pm$  1.77 standard deviation; SD) participated at three test sessions. Participant selection  
149 involved administering a hearing evaluation questionnaire, followed by PTA and tympanometry  
150 during the initial session. Individuals with known hearing disorders, a history of ear surgery, or  
151 tinnitus were excluded. The study encompassed three distinct sessions denoted as session 1, 2,  
152 and 3. Between each consecutive session, a time interval of two to three days was maintained,  
153 with the exception of two participants who had a 14- and 15-day interval between session 2 and  
154 3. Throughout these intervals, participants were instructed to abstain from exposure to loud  
155 activities. During the first session, the better-performing ear was selected based on otoscopy,  
156 tympanometry and PTA. During each session, participants completed a comprehensive test

157 battery consisting of (EHF) PTA, speech in quiet (SPIQ)- and speech in noise (SPiN)-tests,  
158 DPOAEs and AEP-measurements. The selection for the right ear was made for 10 participants,  
159 while the left ear was tested for five participants. As part of Covid-19 safety measures, subjects  
160 wore a face mask during the measurements. The test protocol had a maximum duration of three  
161 hours, and tests were administered in a consistent sequence for all subjects across all sessions.  
162 This study received approval from the UZ Gent ethical committee (BC-05214) and adhered to  
163 the ethical principles outlined in the Declaration of Helsinki. All participants were informed about  
164 the testing procedures, provided informed consent, and received a financial compensation for  
165 their involvement.

#### 166 OTOSCOPY AND TYMPANOMETRY

167 Otoscopy of the ear canal and the tympanic membrane was performed using a Heine Beta 200  
168 LED otoscope (Dover, USA), and showed bilateral normal otoscopic in all subjects. Middle-ear  
169 admittance was bilaterally measured, followed by unilateral measurements (best ear) in the  
170 follow-up sessions, using a GSI TympStar (Grason-Stadler) tympanometer (Minneapolis, USA)  
171 with a 226 Hz, 85 dB sound pressure level (HL) probe tone. All tympanograms were defined as  
172 a type-A according to the Liden-Jerger classification (Jerger, 1970; Lidén, 1969).

#### 173 PURE-TONE AUDIOMETRY

174 Pure-tone thresholds were determined in a double-walled sound-attenuating booth by the use of  
175 an Equinox Interacoustics audiometer (Middelfart, Denmark). Stimuli were transmitted using  
176 Interacoustics TDH-39 headphones (Middelfart, Denmark) and Sennheiser HDA-200  
177 headphones (Wedemark, Germany) for conventional frequencies and EHF, respectively. Air-  
178 conduction thresholds were measured using the modified Hughson-Westlake procedure at  
179 conventional octave frequencies 0.125, 0.250, 0.500, 1, 2, 4 and 8 kHz, half-octave frequencies  
180 3 and 6 kHz, and EHF 10, 12.5, 14, 16, and 20 kHz. Both ears were tested at the first session to

181 select the test-ear for the successive measurements; i.e. ear with better thresholds on conventional  
182 frequencies. All test-ears had thresholds at conventional frequencies of 25 dB HL or better.

### 183 SPEECH INTELLIGIBILITY IN QUIET AND IN NOISE (SPiQ AND SPiN)

184 During each session, SPiQ- and SPiN-tests were conducted within a quiet testing room, using the  
185 Flemish Matrix sentence test (Luts et al. 2014) and Apex 3 software (Francart et al., 2008). The  
186 sentences were presented to the better-performing ear, utilizing a laptop connected to a Fireface  
187 UCX soundcard (RME) (Haimhausen, Germany) and HDA-300 (Sennheiser) headphones  
188 (Wedemark, Germany). In every session, four test lists were randomly selected and presented in  
189 an arbitrary sequence. To mitigate the potential learning effect, two training lists were presented  
190 in the BB-noise-condition (Luts et al., 2014). The test lists encompassed four distinct conditions;  
191 broadband (BB) as well as high-pass (HP) filtered conditions, both in quiet and in noise.  
192 However, due to protocol modifications, the BB condition in quiet was not executed for two  
193 participants. For the BB speech in quiet parameter, no filtering was applied to the speech signal.  
194 However, for the HP-filtered speech in quiet parameter, the speech signal was filtered using a  
195 zero-phase 1024th-order FIR high-pass (HP) filter with a cutoff value of 1650 Hz.

196 The speech in noise parameters were evaluated within the presence of speech-shaped noise  
197 maintained at a constant level of 70 dB SPL. Two distinct noise parameters were employed: the  
198 HP-noise parameter, involving the filtration of both speech and noise signals using the same  
199 cutoff values as the HP-quiet parameter, and the BB-noise parameter, wherein no filtering was  
200 applied to either the speech or noise signals.

201 The matrix-test encompassed a corpus of 50 words, categorized into 10 names, 10 verbs, 10  
202 numerals, 10 adjectives, and 10 nouns. All sentences shared identical syntactical structures, and  
203 the semantic content remained unpredictable. The adaptive procedure outlined by Brand &  
204 Kollmeier, 2002 was used for all test lists, implementing a staircase paradigm to ascertain the



205 speech-reception threshold. The speech level was adjusted by a maximum of 5 dB, progressively  
206 decreasing to a minimal step size of 0.1 dB as the test progressed. The noise was maintained at a  
207 fixed level of 70 dB SPL, with the procedure commencing at a SNR of -20 dB and -4 dB for the  
208 speech-in-quiet and -noise conditions, respectively.

209 For all tests lists, subjects were instructed to repeat the five-word sentences in a forced-choice  
210 setting, providing 10 options for each word. The mean signal level or mean SNR from the six last  
211 reversals was utilized to determine the speech reception threshold (SRT) for the SPiQ- and SPiN-  
212 tests, respectively, wherein lower SRT-values indicated better speech intelligibility.

### 213 DISTORTION PRODUCT OTO-ACOUSTIC EMISSIONS (DPOAEs)

214 During each session, DPOAE measurements were carried out on the designated ear,  
215 encompassing DP-grams and input-output functions. DPOAEs were collected in a quiet testing  
216 room, employing the Universal Smart Box (Intelligent Hearing Systems IHS) (Miami, United  
217 States). To ensure controlled conditions, both ears were shielded using earmuffs (Busters)  
218 (Kontich, Belgium) that were placed on top of a 10D IHS OAE-probe (Miami, United States).  
219 DPOAE responses and noise amplitudes were quantified using the simultaneous presentation of  
220 two primary tones, with f1 and f2 featuring a frequency ratio f2/f1 of 1.22. Noise artifact rejection  
221 was set at 10 dB SPL, and a total of 32 sweeps were recorded for each frequency or input-output  
222 level.

223 DP-grams were obtained with a primary tone level combination of L1/L2 = 65/65 dB SPL and f2  
224 ranging from 553 to 8837 Hz at two points per octave, and from 8837 to 11459 Hz at eight points  
225 per octave, resulting in twelve frequency bands with center frequencies 0.5, 0.7, 1.0, 1.4, 2.0, 2.8,  
226 4.0, 5.7, 8.0, 8.7, 9.5, and 10.3 kHz. The evaluation of DPOAEs was subdivided into three  
227 evaluation methods using commonly used inclusion criteria, i.e. response amplitude  $\geq$  the noise  
228 floor; response amplitude  $\geq$  2SD above the noise floor; response amplitude  $\geq$  6 dB above the

229 noise floor. When responses did not meet the inclusion criteria, the amplitudes were set to the  
230 noise floor levels.

231 Input-output functions were obtained at octave frequencies between 0.5 to 8 kHz (i.e.  $\sqrt{f_1 \cdot f_2} =$   
232 501, 1000, 2000, 3998, 8001, and 10376 Hz) with L2 ranging from 35 to 70 dB SPL in steps of  
233 5 dB. L1/L2 varied across L2 intensities using the scissor paradigm of Kummer et al. (1998)  
234 whereby  $L_1 = 0,4 L_2 + 39$  dB. Extrapolation and non-linear regression were used to estimate DP-  
235 thresholds in which a cubic function was fit to the I/O functions of DPOAE measurements of  
236 each frequency following the method of Verhulst et al. (2016). This way, DPOAE thresholds  
237 were determined as the level of L2 at which the curve reached the distortion component of -25  
238 dB SPL. Thresholds outside the range of -10 – 60 dB, were excluded (Boege & Janssen, 2002),  
239 since these responses are not considered as valid.

#### 240 AUDITORY EVOKED POTENTIAL (AEP) MEASUREMENTS

241 AEP measurements, including EFRs and ABRs were conducted at the test ear using the IHS  
242 universal Smart box and SEPCAM software (Miami, United States). Recordings were performed  
243 in a quiet testing room, with subjects seated in a reclining chair, watching a muted video while  
244 resting their heads on a soft pillow. To minimize alpha-wave interference, subjects were  
245 instructed to relax without falling asleep. Controlled conditions within the hospital setting were  
246 maintained by shielding both ears with earmuffs (Busters) (Kontich, Belgium), turning off  
247 extraneous lights and electronic devices, and applying NuPrep gel for skin preparation.  
248 Disposable Ambu Neuroline electrodes (Ballerup, Denmark) were placed on the vertex (inverting  
249 electrode), nasal flank on the non-test ear side (ground electrode), and bilateral mastoids (non-  
250 inverting electrodes). Electrode impedances were kept below 3 k $\Omega$ , and auditory stimuli were  
251 presented using etymotic ER-2 ear-probes (Chicago, USA).

252 EFRs were evoked using two stimulus types, distinguished by their modulation waveform, i.e.  
253 a sinusoidal amplitude modulated (SAM)-stimulus with a carrier frequency of 4 kHz, and  
254 rectangularly amplitude modulated (RAM)-stimuli, with carrier frequencies 4 and 6 kHz, and a  
255 duty cycle of 25% (Van Der Biest et al., 2023; Vasilkov et al., 2021). EFRs were evoked using  
256 1000 alternating polarity sweeps. Stimuli had a modulation frequency of 110 Hz, a modulation  
257 depth of 100% and a duration of 500 ms which were presented at a rate of 2 Hz. The RAM stimuli  
258 with different carriers were calibrated in such a way to have the same peak-to-peak amplitude as  
259 a 70 dB SPL SAM-tone (carrier: 4 kHz, modulation frequency: 110 Hz, modulation depth:  
260 100%). In this regard, the calibrated RAM stimuli with different carrier frequencies were  
261 presented at 68.24 dB SPL and had the same peak-to-peak amplitudes.

262 The EFR processing was performed in Matlab R2018b. Firstly, the recordings were filtered using  
263 a bandpass filter with low and high cutoff frequencies of 30 Hz and 1500 Hz, respectively. After  
264 filtering the EFRs, epoching and baseline correction was performed. Lastly, a bootstrapping  
265 approach according to Zhu et al. (2013) was adopted in the frequency domain to estimate the  
266 noise-floor and variability of the EFR, as detailed in Keshishzadeh et al. (2020). Subsequently,  
267 EFR-strengths represented the summation of the signal-to-noise spectral magnitude at the  
268 fundamental frequency and its following three harmonics, i.e. 110, 220, 330 and 440 Hz  
269 (Vasilkov et al., 2021).

270 ABRs were evoked using 4000 alternating polarity sweeps of six stimulus types, i.e. three  
271 broadband 80- $\mu$ s clicks presented at levels of 70, 80 and 90 dBpeSPL and three narrowband  
272 toneburst (TB)-stimuli at 0.5 kHz, 1 kHz and 4 kHz with a stimulus duration of 5 ms, 4 ms and  
273 2 ms, respectively. Clicks were presented at a rate of 11 Hz and TBs had a rate of 20 Hz. Data-  
274 processing was performed in Matlab R2018b. ABR recordings were filtered offline between 100  
275 and 1500 Hz using a zero-phase filter. Afterwards, epoching and baseline correction was  
276 performed akin to the method described for EFR processing. After baseline correction, epochs

277 were averaged to yield the ABR waveform. ABR wave I, III and V were manually peak-picked  
278 by an audiologist to identify the respective ABR peak amplitudes ( $\mu\text{V}$ ) and latencies (ms). ABR  
279 amplitudes were defined peak to baseline.

## 280 STATISTICAL ANALYSIS

281 All statistical analyses were conducted using SPSS (IBM) version 25.0. The data-analysis  
282 encompassed a four-tiered approach, comprising one-way repeated measures ANOVA, two-way  
283 random average measures intraclass correlation coefficient (ICC), standard errors of  
284 measurement (SEM), and individual 95% confidence intervals (95%CI). Firstly, a One-way  
285 repeated measures ANOVA was employed to assess variations in PTA, SPiQ and SPiN, DPOAEs  
286 and AEP outcomes across three consecutive measurements. Descriptive parameters were  
287 examined, and the assumptions of the One-way repeated measures ANOVA were verified. When  
288 the significance level ( $p < 0.01$ ) was reached, post-hoc tests were performed to ascertain inter-  
289 session differences. An adjusted significance level of  $p = 0.01$  was selected to protect against  
290 type I errors, due to the large number of variables and statistical analyses performed within the  
291 current study (Gilchrist & Samuels, 2014; Moran, 2003). Secondly, two-way random-average-  
292 measures intraclass correlation coefficients were computed to determine the relative consistency,  
293 i.e. the consistency of the position of individual scores relative to others. The interpretation of  
294 ICC-values followed the classification system proposed by Koo and Li (2016): excellent ICC  
295 ( $>0.90$ ), good ICC (0.75 - 0.90), moderate ICC (0.50 – 0.75) and poor ICC ( $< 0.50$ ). Thirdly,  
296 SEM-scores were calculated to represent the reliability within repeated measures for an  
297 individual subject, reflecting absolute consistency. The latter is calculated as  $\text{SEM} = s \cdot \sqrt{(1-\text{ICC})}$ ,  
298 where 's' represents the standard deviation of all measurements. Finally, given that the substantial  
299 intersubject variability observed in each measure had an influence on the group-based test-retest  
300 95% confidence intervals (CI), we additionally computed 95% CIs of the repeated measures for  
301 each individual separately to visually assess the reliability of different hearing parameters in

302 comparison to each other. This process entailed calculating 95% CIs across measurement  
303 sessions for each parameter and subject. The resulting distribution of individual CIs across  
304 subjects, is visualized using Kernel Density Estimation (KDE) plots, representing the upper and  
305 lower bounds of obtained test-retest CIs for each measure and each subject. These KDE-plots  
306 served to illustrate the variability of test-retest CIs across subjects and enhance the interpretation  
307 of test-retest variations within the data.

## 308 RESULTS

### 309 PURE-TONE AUDIOMETRY (PTA)

310 At session 1, the mean pure-tone average at 3, 4, 6 and 8 kHz (PTA3-8kHz) was 8.08 (SD 6.663,  
311 range -2.00 – 18.00) and the mean pure-tone average at 10, 12.5, 14, 16 and 20 kHz (PTA EHF10-  
312 20 kHz) was 4.40 (SD 8.382, range -8.00 – 18.00). One-way repeated measures ANOVA revealed  
313 no significant changes in pure-tone thresholds between measurements, except for the 0.250 Hz  
314 auditory thresholds [ $F(2, 28) = 6.526, p = 0.005$ ]. Pairwise comparisons indicated a significant  
315 change of 4.00 dB from session 1 to session 3 ( $p = 0.003$ ). Table 1 presents the averages per  
316 session and frequency, as well as the ICCs and SEMs for each tested frequency. In general, good  
317 to excellent ICCs with highly significant between subjects reliability ( $p < 0.001$ ) were obtained  
318 and small SEMs were observed, with the exception of 6, 8 and 20 kHz.

### 319 SPEECH IN QUIET AND SPEECH IN NOISE (SPiQ- AND SPiN)

320 The distribution of speech reception thresholds for both SPiQ-and SPiN-tests across subjects is  
321 depicted in Figures 1A and 1B, respectively. One-way repeated measures ANOVA indicated no  
322 significant alterations in SRT-values between measurement sessions for the BB conditions in  
323 quiet [ $F(2, 24) = 1.549, p > 0.01$ ], nor in noise [ $F(2, 28) = 0.690, p > 0.01$ ]. In contrast, significant  
324 changes in SRT- and SNR-values were found for the HP condition in quiet [ $F(2, 28) = 12.266, p$   
325  $< 0.001$ ], and the HP condition in noise [ $F(2, 28) = 7.788, p = 0.002$ ]. Pairwise comparisons  
326 unveiled SRT-improvements for the HP quiet condition between each session with a significant

327 change of 1.72 dB SPL from session 1 to session 2 ( $p = 0.011$ ), 1.12 dB SPL from session 2 to  
328 session 3 ( $p = 0.013$ ) and 2.84 dB SPL from session 1 to session 3 ( $p = 0.001$ ). Furthermore, a  
329 significant SNR-change of 1.13 dB SNR between session 2 and 3 ( $p < 0.001$ ), and a significant  
330 SNR-change of 1.21 dB SNR between session 1 and 3 ( $p = 0.003$ ) was found for the HP noise  
331 condition. It should be noted that the initial training procedure before the start of the  
332 measurements only incorporated BB-speech in noise.

333 Table 2 displays ICCs and SEMs for SPiQ and SPiN tests. Moderate ICCs with highly  
334 significant between-subject variability ( $p < 0.001$ ) and small SEMs were observed for HP  
335 filtered conditions. In the BB noise-condition, a moderate ICC and small SEM were retained.  
336 However, the BB condition in quiet exhibited a poor ICC of 0.341 and a high SEM.

337 DISTORTION PRODUCT OTOACOUSTIC EMISSIONS (DPOAEs)

338 DP-GRAM

339 Per criterium, i.e.  $SNR \geq 0$ ,  $SNR \geq 2$  SD, and  $SNR \geq 6$  dB, one-way repeated measures  
340 ANOVA indicated no significant changes in DP-gram amplitudes ( $p > 0.01$ ) for all tested  
341 frequencies. ICCs and SEMs are shown in Table 3. Overall, the  $SNR \geq 0$ -criterium showed the  
342 highest ICCs (moderate-to-good), followed by the  $SNR \geq 6$  dB-criterium and the  $SNR \geq 2SD$ -  
343 criterium, respectively. The latter criterium is additionally characterized by greater variability  
344 among the different tested frequencies. Secondly, remarkably worse ICCs were found for the  
345 lower frequencies of 1 and 1.5 kHz. SEMs showed relatively large values overall, with the  
346  $SNR \geq 2SD$ -criterium showing the largest values relative to the other criteria. Figures 2 A, B,  
347 and C depict KDE-plots of the zero-criterium, illustrating the distribution of individual test-  
348 retest 95% CIs for different measures. A sharp peak in the KDE signifies a more concentrated  
349 distribution of test-retest CIs, indicating good overall reliability across the test population.  
350 Conversely, a broader peak implies increased variability of individual test-retest CIs across  
351 individuals, reflecting a lower reliability for the corresponding parameter across the population.

352 INPUT-OUTPUT FUNCTION

353 One-way repeated measures ANOVA indicated no significant changes ( $p > 0.01$ ) for all tested  
354 frequencies, computed per inclusion criterium. ICCs and SEMs are shown in Table 4.  
355 The  $SNR \geq 6\text{dB}$ -criterium showed overall the highest ICCs, followed by  $SNR \geq 2SD$  and  
356  $SNR \geq 0$ , retaining both very similar ICCs and SEMs. However, overall very poor ICCs, in  
357 addition to large SEMs, were observed across the six tested frequencies for DPOAE input-  
358 output measures.

359 DP-GRAMS AND INPUT-OUTPUT FUNCTIONS IN RELATION TO PURE-TONE AUDITORY THRESHOLDS

360 ICCs and SEMs of DP-grams and input output functions, in relation to PTA are illustrated in  
361 Figures 3A, B, C, and D. In terms of DP-gram ICCs (A), a pattern of generally lower but more  
362 consistent outcomes across different frequencies and evaluation criteria was observed  
363 compared to PTA. Notably, exceptions were observed at 6 kHz and, predominantly, 8 kHz,  
364 where DP-grams exhibited better ICCs than the audiogram. This trend corresponded with the  
365 DP-gram SEMs (C). The DPOAE input-output functions exhibited notably lower ICCs when  
366 juxtaposed with pure-tone audiometry (B). Moreover, increased variability across different  
367 evaluation criteria was evident, particularly with improved outcomes for the 6 dB criterion.  
368 SEMs confirmed less favorable results for DP-thresholds (D), emphasizing the 6 dB criterion  
369 as the most reliable.

370 AUDITORY EVOKED POTENTIALS (AEP)

371 AUDITORY BRAINSTEM RESPONSES (ABR)

372 A one-way repeated measures ANOVA revealed no significant differences in click-ABR  
373 amplitudes and latencies of peaks I, III, and V at 70 dBpeSPL, 80 dBpeSPL, and 90 dBpeSPL,  
374 as well as for TB-stimuli, across measurement sessions ( $p > 0.01$ ). Tables 5 and 6 display ICCs  
375 and SEMs for click- and TB-stimuli peak I, III, and V amplitudes, and latencies. Generally, good-  
376 to-excellent ICCs with highly significant between-subject variances and small SEMs were  
377 observed for click-ABR peak I-, III- and V-latencies. These findings align with click-ABR peak

378 V-amplitudes, showing good ICCs with highly significant between-subject variances, except for  
379 the click-ABR 70 dBpeSPL, retaining a moderate ICC. In contrast to peaks I- and III-click  
380 latencies, peak I- and III-amplitudes showed moderate-to-very poor ICCs, with poor average  
381 measures.

382 TB-amplitudes generally exhibited slightly lower ICCs and higher SEMs compared to clicks.  
383 Similar to click-stimuli findings, peak I amplitudes exhibited moderate-to-very poor ICCs. For  
384 latencies, moderate-to-good ICCs with highly significant between-subject variances ( $p < 0.001$ )  
385 and small SEMs were observed for peaks I- and III- and V-latencies. Figures 2 D, E, F and G  
386 display KDE-plots of ABR-amplitudes, and -latencies, representing the distribution of individual  
387 test-retest 95% CIs for different measures.

#### 388 ENVELOPE FOLLOWING RESPONSES (EFR)

389 Figure 1C depicts the EFR-strength distribution for SAM- and RAM-stimuli across the three  
390 consecutive sessions. Consistent with prior findings (Vasilkov et al., 2021), EFRs were stronger  
391 for the RAM stimuli than the SAM stimulus. An outlier identified in the SAM-evoked response  
392 during session two was excluded due to potential data corruption caused by a 50 Hz noise from  
393 nearby electrical equipment. No significant changes in EFR-strength were found between  
394 measurements for the SAM-stimulus [ $F(2, 26) = 0.066, p > 0.01$ ], and RAM-stimuli; i.e. RAM 4  
395 kHz [ $F(2, 28) = 0.383, p > 0.01$ ] and RAM 6 kHz [ $F(2, 28) = 1.299, p > 0.01$ ]. Moreover, the  
396 ICC demonstrated excellent test-retest reliability for both the SAM and RAM stimuli at 4 kHz  
397 and 6 kHz, yielding average measures of 0.882 (95% BI [0.708;0.959];  $F(13, 26) = 7.975, p <$   
398  $0.001$ ), 0.950 (95% BI [0.883;0.982];  $F(14, 28) = 19.355, p < 0.001$ ) and 0.930 (95% BI  
399 [0.837;0.974];  $F(14, 28) = 14.553, p < 0.001$ ), respectively. Table 7 provides detailed information  
400 on ICCs and SEMs.

401 Figure 4 depicts (non)-significant individual EFR-changes across the three test sessions, where  
402 the 95% confidence intervals (CI) of the individual EFR-strength differences are computed using



403 all EFR stimulus averages from the two respective sessions through a bootstrapping method.  
404 Individual EFR-strengths are considered significantly different between two sessions if their CI  
405 does not overlap with the zero-line. Gray dashed lines represent the overall CI, calculated as the  
406 mean CI across all subjects and measurements, indicating significance when individual  
407 datapoints fall outside this interval. Notably, only subject seven exhibited significant alterations  
408 in both-RAM-evoked EFR-amplitudes. It is worth mentioning that subject 7 showed signs of  
409 agitation toward the end of the comprehensive protocol, particularly in response to RAM-stimuli.  
410 This observation appeared to correspond with a less stable EEG-signal, as recorded in the  
411 logbook.

412 To provide a more comprehensive perspective on observed variations in ABR- and EFR-  
413 magnitudes, both considered as potential EEG-markers of CS, Figure 5 displays individual  
414 strengths and distribution boxplots. Additionally, Figure 2H displays KDE-plots for both  
415 parameters, illustrating the distribution of individual 95% CIs calculated across three sessions.  
416 Both figures highlight the superior reliability of EFR-magnitudes compared to ABR-  
417 amplitudes.

## 418 DISCUSSION

### 419 HIGH RELIABILITY OF PURE-TONE AUDIOMETRY WITH CONSIDERATION FOR FREQUENCY SPECIFIC VARIATIONS

420 Prior research has recommended using a frequency range up to 14 kHz for monitoring purposes  
421 (Rodríguez Valiente et al., 2014), as frequencies beyond this threshold show substantial intra-  
422 subject threshold variability (Frank, 2001; Schmuziger et al., 2007). While the current study  
423 revealed good-to-excellent ICCs for both conventional and extended high frequencies,  
424 frequency-specific disparities should be taken into consideration within clinical practice since  
425 6, 8, and 20 kHz retained moderate-to-poor ICCs. The higher variability at 6 and 8 kHz aligns  
426 with the findings of Schlauch and Carney (2011), and may be linked to suboptimal earphone  
427 positioning or calibration methods. The increased variability at 20 kHz is likely due to standing

428 waves, suggesting that extending the frequency range up to 20 kHz is not advisable for  
429 monitoring purposes.

430 The good test-retest reliability of (high-frequency) audiometry is corroborated by studies  
431 conducted by Swanepoel et al. (2010) and Ishak et al. (2011), as well as by several other  
432 investigations that employed diverse transducer models (Fausti et al., 1998; Frank, 1990, 2001;  
433 Frank & Dreisbach, 1991; Schmuziger et al., 2004).

434 Repeated-measures ANOVA revealed significant differences in 250 Hz thresholds. The need  
435 for subjects to wear masks, as part of Covid-19 safety measures, may have affected response  
436 accuracy, especially in lower frequency tests. The observed improvement in results during  
437 session three may be attributed to subjects becoming more accustomed to the potential masking  
438 effect caused by wearing masks.

439 NEED FOR A RELIABLE SPiQ-AND SPiN TEST.

440 SPiQ- and SPiN-tests showed no significant differences between sessions for the BB-conditions,  
441 while significant SRT- and SNR-changes were found for the HP-conditions in quiet and in noise.  
442 This learning effect was previously documented by Luts et al. (2014), highlighting a large  
443 decrease in SRT occurring between the first and the second measurement which decreased to a  
444 value below 1 dB after the second list. Similar trends were observed for all language-specific  
445 tests covered in the review paper by Kollmeier et al. (2015), suggesting that the training effect  
446 might be associated more with the nature of the task and test structure rather than language-  
447 specific characteristics. However, in the current study, a training effect was observed despite the  
448 provision of two training lists. Firstly, the inclusion of two HP filtered training lists might have  
449 counteracted learning, as only BB training lists were intended. However, presenting multiple  
450 training lists extends the test duration, affecting subjects' attention span and potentially  
451 influencing outcomes. Nevertheless, Maele et al. (2021) reported significant SNR-improvements

452 in all tested conditions, including BB, even with two BB-training lists provided. Secondly, a  
453 closed-set test format might contribute to the learning effect, as subjects could more easily learn  
454 words when both heard and visualized. The intention of displaying the possible words was to  
455 mitigate potential performances improvement across sessions, as subjects are aware of the words  
456 they may encounter from the outset. Nonetheless, the review paper of Kollmeier et al. (2015)  
457 reported a training effect in both open- and closed-set test formats for each language examined.  
458 Thirdly, within this study, participants were directed to respond in a forced-choice format to  
459 mitigate the learning effect. This was prompted by the observation that during the initial session,  
460 subjects frequently signaled non-detection of the word more swiftly, yet exhibited increasing  
461 confidence in subsequent measurements. This behavior might lead to speculation and potentially  
462 improved performance in subsequent sessions, which might indirectly contribute to better results.  
463 In sum, while further investigation is needed, caution is advised in using the matrix test in  
464 repeated measures or monitoring, due to the potential influence of learning effects.

465 DP-GRAMS SHOW GREAT VARIABILITY.

466 One-way-repeated measures ANOVA indicated no significant changes in DP-amplitudes  
467 between the measurements for all tested criteria and frequencies. These findings align with a  
468 study conducted in 2010, reporting no significant differences within time-intervals up to 60  
469 minutes (Keppler et al., 2010). However, when the time-interval extended to 7 days, a  
470 significant difference was noted, suggesting decreased reliability of DP-grams with increased  
471 time intervals (Keppler et al., 2010). Engdahl et al. (1994) and Wagner et al. (2008) noted that  
472 prolonged time intervals lead to increased standard deviations due to greater variation in middle  
473 ear pressure, room- and biological noise. Probe refitting at each session on different days  
474 further contributes to variability, as indicated by research highlighting the impact of probe  
475 replacement on the level of background noise and acoustic leakage (Beattie & Bleech, 2000;

476 Beattie et al., 2003; Franklin et al., 1992; Keppler et al., 2010; Mills et al., 2007; Wagner et al.,  
477 2008; Zhao & Stephens, 1999).

478 SEM values in this study generally surpassed those reported in other studies (Beattie et al.,  
479 2003; Franklin et al., 1992; Keppler et al., 2010; Ng & Mcpherson, 2005; Wagner et al., 2008).

480 Various factors could account for this increased variability. Firstly, subject-generated noise has  
481 the potential to impact DP-gram amplitudes, introducing variability due to differences in  
482 patient cooperation and ear canal acoustics (Keppler et al., 2010; Wagner et al., 2008).

483 Secondly, equipment-related noise (Keppler et al., 2010) and recording parameters could  
484 significantly influence DP-gram amplitudes. Franklin et al. (1992) demonstrated lower  
485 reliability for DPOAE amplitudes elicited using lower primary tone level combinations (L1/L2

486 = 65/55 dB SPL), as used in this study, compared to higher primary level combinations (L1/L2  
487 = 75/70 dB SPL). Another study in 1996, evaluating two stimulus protocols, consistently  
488 reported larger absolute DP-amplitudes for higher primary intensities, with greater variability

489 for lower primary level combinations relative to higher intensities (Hall, 2000). Thirdly,  
490 analysis strategies vary significantly across studies. In contrast to previous research, this study  
491 did not exclude responses that did not meet inclusion criteria. Instead, amplitudes were adjusted

492 to the minimum level (i.e. the noise floor level), resulting in larger standard deviations, and  
493 consequently, larger SEMs. Additionally, the reliability varies with different inclusion criteria,  
494 as discussed earlier, and across frequencies, as depicted in Table 3. The notably higher standard

495 deviations observed at frequencies 1.0 and 1.5 kHz, compared to others, are likely attributed to  
496 low-frequency noise contaminations (Beattie et al., 2003; Keppler et al., 2010; Wagner et al.,  
497 2008; Zhao & Stephens, 1999). The increased variability in DP-gram amplitudes at higher

498 frequencies is probably caused by ear-canal acoustics, particularly standing waves, amplifying  
499 intrinsic variability due to differences in sound pressure at the tympanic membrane and probe  
500 microphone (Keppler et al., 2010; Mills et al., 2007).

501 CLINICAL DPOAE INPUT-OUTPUT FUNCTIONS ARE NOT RELIABLE.

502 Previous research has indicated the potential of estimated DPOAE thresholds to predict pure-tone  
503 thresholds (Boege & Janssen, 2002; Goldman et al., 2006; Gorga et al., 2003). However, studies  
504 on test-retest reliability are limited. In the current study, one-way repeated measures ANOVA  
505 revealed no significant changes in DP-thresholds between measurements for all tested criteria  
506 and frequencies. Nevertheless, ICCs and SEMs yielded very poor results, suggesting a low level  
507 of reliability. The highly variable nature of the I/O function among subjects and even for different  
508 stimuli at different frequencies, as highlighted by Kimberley and Nelson (1989), Hall (2000), and  
509 Harris (1990), questions the clinical utility of this method. Additionally, Harris (1990)  
510 emphasized the need for strict minimum noise requirements for reliable responses,  
511 recommending measures such as conducting DPOAE measurements in a sound-attenuating  
512 booth, setting test protocol stopping criteria for a very low noise level, and employing continuous  
513 signal averaging until the minimum noise level is reached. Popelka et al. (1993) reported that  
514 achieving a noise floor of -40 dB for recording a single I/O function may take up to 45 minutes  
515 of testing. In summary, both the current study and the findings of Popelka et al. (1993) suggest  
516 that, in its present form, input-output functions may not be considered useful in clinical settings.

517 EFR MEASURES YIELD BETTER TEST-RETEST RELIABILITY, RELATIVE TO THE ABR.

518 Although it has been reported that the ABR is stable over long time periods, a large number of  
519 studies refers to peak V, which is often a more reliable and robust parameter relative to peak I;  
520 characterized by smaller amplitudes and hence increased variability across sessions (Prendergast  
521 et al., 2018). The reliability of wave I is however of great value, as wave I-amplitude has been  
522 identified as a potential non-invasive measure of CS (Mehraei et al., 2016).

523 The present study showed no significant changes for peaks I, III and V between the three different  
524 test sessions. These results align with a 2018 study, showing no significant changes in click- and  
525 speech-evoked brainstem responses across test sessions (Bidelman et al., 2018). Additionally,

526 Munjal et al. (2016), evaluating the reliability of the absolute latency of waves I, III and V and  
527 interpeak latencies, showed good test-retest reliability for all response parameters, except for the  
528 absolute latency of wave I.

529 The analyses of ICCs and SEMs unveiled several trends among the different waves and stimuli.  
530 Firstly, a higher within-subject reliability was found for peak V relative to peak I. These results  
531 are in line with the study of Lauter and Karzon (1990), who reported low level of consistency  
532 across subjects for peak I of ABR. Sininger and Cone-Wesson (2002) have shown that peripheral  
533 hearing and testing parameters, amongst others ambient noise and minimal wax in the external  
534 auditory canal, can affect the latency of wave I in ABR-measurements. Secondly, peak V is  
535 assumed to be more robust because of its greater amplitude compared to peak I. And thirdly, peak  
536 I may potentially be reduced due to CS and OHC-damage, while peak V may be enhanced due  
537 to central gain mechanisms (Auerbach et al., 2014; Gu et al., 2012; Schaette & McAlpine, 2011).  
538 Therefore, the interpretation of wave I-amplitudes and latencies requires some caution.

539 In addition to greater reliability for peak V relative to peak I, the present study also retained  
540 smaller ICCs and larger SEMs for click-amplitudes compared to click-latencies, consistent with  
541 Bidelman et al. (2018). Negligibly small intra-subject variability in ABR latencies are in addition  
542 in agreement with previous studies of Edwards et al. (1982) and Oyler et al. (1991). Firstly, these  
543 results could be attributed to the greater magnitudes of latencies compared to amplitudes,  
544 contributing to enhanced robustness in test outcomes. Secondly, evoked potential amplitudes are  
545 susceptible to nonbiological factors, such as electrode impedance and orientation relative to  
546 source generators. This suggest that the amplitude might be a poor metric for reliably assessing  
547 subtle changes in ABR-measurements with certain experimental manipulations, including noise  
548 exposure, ototoxicity, age, and training (Bidelman et al., 2018). The use of ear canal tiptrodes, as  
549 opposed to scalp mounted electrodes could result in higher reliability since the recording site has  
550 moved closer to the generator of wave I, specifically the auditory nerve. This assumption aligns

551 with the study of Bauch and Olsen (1990), showing increased wave I-amplitudes with ear canal  
552 tiptrodes compared to mastoid electrodes, and Bieber et al. (2020), reporting good-to-excellent  
553 wave I and wave V amplitude ICCs when measured from the ear canal. However, the study of  
554 Prendergast et al. (2018) demonstrated only a small increase in reliability for waves I and V when  
555 using canal tiptrodes compared to mastoid electrodes. The benefits for the summation potential  
556 however, were greater. In sum, although wave I has proved valuable, particularly in research  
557 studies, as a more direct measure of peripheral auditory function, low amplitude ICCs were found  
558 in the present study, questioning the applicability of clinical waveform interpretation when  
559 recorded under aforementioned conditions. Moreover, prior studies showing good-to-excellent  
560 test-retest reliability did not address clinical feasibility either, due to the extensive test duration,  
561 involving 5600 to 10000 sweeps (Bieber et al., 2020; Guest et al., 2019; Prendergast et al., 2018).

562 When comparing broadband clicks with toneburst-stimuli, clicks generally retained bigger  
563 responses and are hence slightly more reliable biomarkers. It is assumed that a broadband click  
564 stimulus drives more fibers simultaneously, potentially eliciting larger and more robust  
565 responses. Subsequently, TBs are identified as less clearly detectable peaks, indicating higher  
566 interrater variability and therefore lower reliability in general. This hypothesis is supported by  
567 generally smaller amplitudes and larger latencies in TB-responses, likely stemming from  
568 narrowed basilar-membrane stimulation (Gorga et al., 1988; Rasetshwane et al., 2013). The  
569 impact of prolonged latencies and decreased amplitudes with narrower BM stimulation is further  
570 pronounced in lower frequencies relative to higher frequencies (500 Hz vs 4 kHz), mainly due to  
571 cochlear wave dispersion (Rasetshwane et al., 2013).

572 As noise-induced CS primarily targets AN fibers with high thresholds, and phase locking to  
573 temporal envelopes is in addition particularly strong in these fibers, the EFR-strength could  
574 potentially be a more robust measure, relative to ABR-amplitudes (Vasilkov et al., 2021).  
575 Additionally, phase information can be extracted from EFRs, and measures of phase-locking

576 values might be less susceptible to anatomical variations in human (Gorga et al., 1988),  
577 generally interfering amplitude measures. These findings align with the results of Bidelman et  
578 al. (2018), showing FFRs that yield overall higher test-retest reliability, relative to their  
579 conventional click-ABR counterparts, and are in harmony with the current study, revealing  
580 good ICCs and small SEMs.

## 581 CONCLUSION

582 The quest to identify noninvasive early markers of noise-induced SNHL in humans has generated  
583 various measures of interest. Nevertheless, comprehensive studies assessing the test-retest  
584 reliability of multiple measures and stimuli within a single study remain limited, and a  
585 standardized clinical protocol encompassing robust noninvasive early markers of SNHL has not  
586 yet been established. In light of these gaps, the present study aimed to explore the intra-subject  
587 variability of various potential noninvasive EEG-biomarkers of CS and other early indicators of  
588 SNHL within the same individuals. The study underscores the need for caution when interpreting  
589 presumed noninvasive SNHL measures. While pure-tone audiometry generally exhibited high  
590 test reliability, frequency-specific differences should be taken into account within clinical  
591 practice. Furthermore, extending the frequency range beyond 16 kHz is not advisable for clinical  
592 use. The observed learning effect in the speech-sentence test highlights the need for prudence  
593 when employing the matrix sentence test in repeated measurements. The variability observed in  
594 DPOAEs necessitates consistent ear probe replacement, meticulous measurement techniques, and  
595 optimal testing conditions to minimize variability in DP-grams, and renders I/O-functions  
596 unsuitable for clinical application. In terms of auditory evoked potentials, EFRs demonstrated  
597 greater reliability compared to ABRs. However, it is crucial to exercise caution due to factors  
598 intrinsic and external to the individual, which may contribute to increased variability across  
599 longitudinal measurements.



## 600 REFERENCES

- 601 Auerbach, B. D., Rodrigues, P. V., & Salvi, R. J. (2014). Central gain control in tinnitus and  
602 hyperacusis. *Frontiers in Neurology*, *5*, 206.
- 603 Bauch, C. D., & Olsen, W. O. (1990). Comparison of ABR amplitudes with TIPtrode™ and mastoid  
604 electrodes. *Ear and Hearing*, *11*(6), 463-467.
- 605 Beattie, R., & Bleech, J. (2000). Effects of sample size on the reliability of noise floor and DPOAE.  
606 *British Journal of Audiology*, *34*(5), 305-309.
- 607 Beattie, R. C., Kenworthy, O., & Luna, C. A. (2003). Immediate and short-term reliability of distortion-  
608 product otoacoustic emissions: Confiabilidad inmediata ya corto plazo de las emisiones  
609 otoacústicas por productos de distorsión. *International Journal of Audiology*, *42*(6), 348-354.
- 610 Bharadwaj, H. M., Masud, S., Mehraei, G., Verhulst, S., & Shinn-Cunningham, B. G. (2015). Individual  
611 differences reveal correlates of hidden hearing deficits. *The Journal of neuroscience : the*  
612 *official journal of the Society for Neuroscience*, *35*(5), 2161-2172.  
613 <https://doi.org/10.1523/JNEUROSCI.3915-14.2015>
- 614 Bharadwaj, H. M., Verhulst, S., Shaheen, L., Liberman, M. C., & Shinn-Cunningham, B. G. (2014).  
615 Cochlear neuropathy and the coding of supra-threshold sound. *Frontiers in Systems*  
616 *Neuroscience*, *8*, 26. <https://doi.org/10.3389/fnsys.2014.00026>
- 617 Bidelman, G. M., Pousson, M., Dugas, C., & Fehrenbach, A. (2018). Test-retest reliability of dual-  
618 recorded brainstem versus cortical auditory-evoked potentials to speech. *Journal of the*  
619 *American Academy of Audiology*, *29*(02), 164-174.
- 620 Bieber, R. E., Fernandez, K., Zalewski, C., Cheng, H., & Brewer, C. C. (2020). Stability of early auditory  
621 evoked potential components over extended test-retest intervals in young adults. *Ear and*  
622 *Hearing*, *41*(6), 1461.
- 623 Boege, P., & Janssen, T. (2002). Pure-tone threshold estimation from extrapolated distortion product  
624 otoacoustic emission I/O-functions in normal and cochlear hearing loss ears. *The Journal of*  
625 *the Acoustical Society of America*, *111*(4), 1810-1818.  
626 <https://www.ncbi.nlm.nih.gov/pubmed/12002865>
- 627 Bourien, J., Tang, Y., Batrel, C., Huet, A., Lenoir, M., Ladrech, S., Desmadryl, G., Nouvian, R., Puel, J.-  
628 L., & Wang, J. (2014). Contribution of auditory nerve fibers to compound action potential of  
629 the auditory nerve. *Journal of Neurophysiology*, *112*(5), 1025-1039.
- 630 Coradini, P. P., Cigana, L., Selistre, S. G., Rosito, L. S., & Brunetto, A. L. (2007). Ototoxicity from  
631 cisplatin therapy in childhood cancer. *Journal of Pediatric Hematology/Oncology*, *29*(6), 355-  
632 360.
- 633 Cruickshanks, K. J., Nondahl, D. M., Tweed, T. S., Wiley, T. L., Klein, B. E., Klein, R., Chappell, R.,  
634 Dalton, D. S., & Nash, S. D. (2010). Education, occupation, noise exposure history and the 10-  
635 yr cumulative incidence of hearing impairment in older adults. *Hearing Research*, *264*(1-2),  
636 3-9.
- 637 Don, M., Ponton, C. W., Eggermont, J. J., & Masuda, A. (1994). Auditory brainstem response (ABR)  
638 peak amplitude variability reflects individual differences in cochlear response times. *The*  
639 *Journal of the Acoustical Society of America*, *96*(6), 3476-3491.
- 640 Edwards, R. M., Buchwald, J. S., Tanguay, P. E., & Schwafel, J. A. (1982). Sources of variability in  
641 auditory brain stem evoked potential measures over time. *Electroencephalography and*  
642 *Clinical Neurophysiology*, *53*(2), 125-132.
- 643 Engdahl, B., Arnesen, A. R., & Mair, I. W. (1994). Reproducibility and short-term variability of  
644 transient evoked otoacoustic emissions. *Scandinavian Audiology*, *23*(2), 99-104.
- 645 Fausti, S. A., Henry, J. A., Hayden, D., Phillips, D. S., & Frey, R. H. (1998). Intrasubject reliability of  
646 high-frequency (9-14 kHz) thresholds: tested separately vs. following conventional-  
647 frequency testing. *Journal of the American Academy of Audiology*, *9*(2).

- 648 Francart, T., van Wieringen, A., & Wouters, J. (2008). APEX 3: a multi-purpose test platform for  
649 auditory psychophysical experiments. *Journal of Neuroscience Methods*, 172(2), 283-293.  
650 <https://doi.org/10.1016/j.jneumeth.2008.04.020>
- 651 Frank, T. (1990). High-frequency hearing thresholds in young adults using a commercially available  
652 audiometer. *Ear and Hearing*, 11(6), 450-454.
- 653 Frank, T. (2001). High-frequency (8 to 16 kHz) reference thresholds and intrasubject threshold  
654 variability relative to ototoxicity criteria using a Sennheiser HDA 200 earphone. *Ear and*  
655 *Hearing*, 22(2), 161-168.
- 656 Frank, T., & Dreisbach, L. E. (1991). Repeatability of high-frequency thresholds. *Ear and Hearing*,  
657 12(4), 294-295.
- 658 Franklin, D. J., McCoy, M. J., Martin, G. K., & Lonsbury-Martin, B. L. (1992). Test/retest reliability of  
659 distortion-product and transiently evoked otoacoustic emissions. *Ear and Hearing*, 13(6),  
660 417-429.
- 661 Furman, A. C., Kujawa, S. G., & Liberman, M. C. (2013). Noise-induced cochlear neuropathy is  
662 selective for fibers with low spontaneous rates. *Journal of Neurophysiology*, 110(3), 577-586.  
663 <https://doi.org/10.1152/jn.00164.2013>
- 664 Gilchrist, M., & Samuels, P. (2014). Statistical hypothesis testing. *Birmingham: Birmingham City*  
665 *University*.
- 666 Goldman, B., Sheppard, L., Kujawa, S. G., & Seixas, N. S. (2006). Modeling distortion product  
667 otoacoustic emission input/output functions using segmented regression. *The Journal of the*  
668 *Acoustical Society of America*, 120(5), 2764-2776.
- 669 Gorga, M. P., Kaminski, J. R., Beauchaine, K. A., & Jesteadt, W. (1988). Auditory brainstem responses  
670 to tone bursts in normally hearing subjects. *Journal of Speech, Language, and Hearing*  
671 *Research*, 31(1), 87-97.
- 672 Gorga, M. P., Neely, S. T., Dorn, P. A., & Hoover, B. M. (2003). Further efforts to predict pure-tone  
673 thresholds from distortion product otoacoustic emission input/output functions. *The Journal*  
674 *of the Acoustical Society of America*, 113(6), 3275-3284.
- 675 Grinn, S. K., Wiseman, K. B., Baker, J. A., & Le Prell, C. G. (2017). Hidden Hearing Loss? No Effect of  
676 Common Recreational Noise Exposure on Cochlear Nerve Response Amplitude in Humans.  
677 *Frontiers in Neuroscience*, 11, 465. <https://doi.org/10.3389/fnins.2017.00465>
- 678 Gu, J. W., Herrmann, B. S., Levine, R. A., & Melcher, J. R. (2012). Brainstem auditory evoked  
679 potentials suggest a role for the ventral cochlear nucleus in tinnitus. *Journal of the*  
680 *Association for Research in Otolaryngology*, 13, 819-833.
- 681 Guest, H., Munro, K. J., & Plack, C. J. (2018). Acoustic Middle-Ear-Muscle-Reflex Thresholds in  
682 Humans with Normal Audiograms: No Relations to Tinnitus, Speech Perception in Noise, or  
683 Noise Exposure. *Neuroscience*. <https://doi.org/10.1016/j.neuroscience.2018.12.019>
- 684 Guest, H., Munro, K. J., Prendergast, G., Howe, S., & Plack, C. J. (2017). Tinnitus with a normal  
685 audiogram: Relation to noise exposure but no evidence for cochlear synaptopathy [Article].  
686 *Hearing Research*, 344, 265-274. <https://doi.org/10.1016/j.heares.2016.12.002>
- 687 Guest, H., Munro, K. J., Prendergast, G., & Plack, C. J. (2019). Reliability and interrelations of seven  
688 proxy measures of cochlear synaptopathy [Article in Press]. *Hearing Research*.  
689 <https://doi.org/10.1016/j.heares.2019.01.018>
- 690 Hall, J. (2000). Distortion product and transient evoked OAEs: Nonpathologic factors influencing  
691 measurement. Handbook of Otoacoustic Emissions. In: Singular Publishing Group. Cengage  
692 Learning, San Diego.
- 693 Harris, F. (1990). Distortion-product otoacoustic emissions in humans with high frequency  
694 sensorineural hearing loss. *Journal of Speech, Language, and Hearing Research*, 33(3), 594-  
695 600.
- 696 Ishak, W. S., Zhao, F., Stephens, D., Culling, J., Bai, Z., & Meyer-Bisch, C. (2011). Test-retest reliability  
697 and validity of Audioscan and Békésy compared with pure tone audiometry. *Audiological*  
698 *Medicine*, 9(1), 40-46.

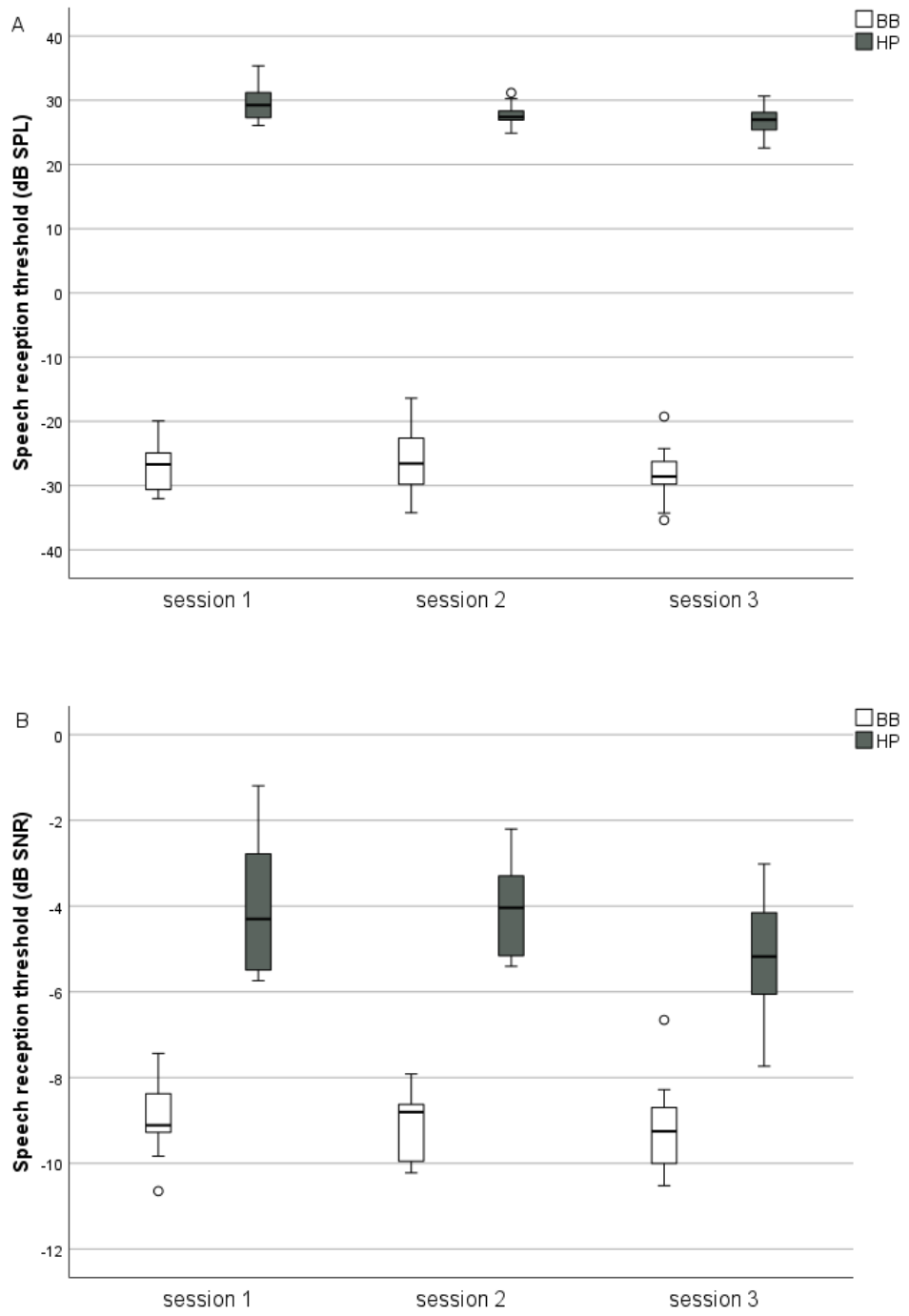
- 699 Jansen, E., Helleman, H., Dreschler, W., & De Laat, J. (2009). Noise induced hearing loss and other  
700 hearing complaints among musicians of symphony orchestras. *International Archives of*  
701 *Occupational and Environmental Health*, 82, 153-164.
- 702 Jerger, J. (1970). Clinical experience with impedance audiometry. *Archives of Otolaryngology*, 92(4),  
703 311-324.
- 704 Keppler, H., Dhooge, I., Maes, L., D'Haenens, W., Bockstael, A., Philips, B., Swinnen, F., & Vinck, B.  
705 (2010). Transient-evoked and distortion product otoacoustic emissions: A short-term test-  
706 retest reliability study. *International Journal of Audiology*, 49(2), 99-109.  
707 <https://doi.org/10.3109/14992020903300431>
- 708 Keshishzadeh, S., Garrett, M., Vasilkov, V., & Verhulst, S. (2020). The derived-band envelope  
709 following response and its sensitivity to sensorineural hearing deficits. *Hearing Research*,  
710 392, 107979.
- 711 Knight, K. R., Kraemer, D. F., Winter, C., & Neuwelt, E. A. (2007). Early changes in auditory function as  
712 a result of platinum chemotherapy: use of extended high-frequency audiometry and evoked  
713 distortion product otoacoustic emissions. *Journal of Clinical Oncology*, 25(10), 1190-1195.
- 714 Kollmeier, B., Warzybok, A., Hochmuth, S., Zokoll, M. A., Uslar, V., Brand, T., & Wagener, K. C. (2015).  
715 The multilingual matrix test: Principles, applications, and comparison across languages: A  
716 review. *International Journal of Audiology*, 54 Suppl 2, 3-16.  
717 <https://doi.org/10.3109/14992027.2015.1020971>
- 718 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation coefficients  
719 for reliability research. *Journal of Chiropractic Medicine*, 15(2), 155-163.
- 720 Kujawa, S. G., & Liberman, M. C. (2009). Adding insult to injury: cochlear nerve degeneration after  
721 "temporary" noise-induced hearing loss. *The Journal of neuroscience : the official journal of*  
722 *the Society for Neuroscience*, 29(45), 14077-14085.  
723 <https://doi.org/10.1523/JNEUROSCI.2845-09.2009>
- 724 Kummer, P., Janssen, T., & Arnold, W. (1998). The level and growth behavior of the 2 f1-f2 distortion  
725 product otoacoustic emission and its relationship to auditory sensitivity in normal hearing  
726 and cochlear hearing loss. *The Journal of the Acoustical Society of America*, 103(6), 3431-  
727 3444.
- 728 Lauter, J. L., & Karzon, R. G. (1990). Individual Differences in Auditory Electric Responses:  
729 Comparisons of Between-Subject and Within-Subject Variability IV. Latency-variability  
730 Comparisons in Early, Middle, and Late Responses. *Scandinavian Audiology*, 19(3), 175-182.
- 731 Liberman, M. C., Epstein, M. J., Cleveland, S. S., Wang, H. B., & Maison, S. F. (2016). Toward a  
732 Differential Diagnosis of Hidden Hearing Loss in Humans [Article]. *PloS One*, 11(9), 15, Article  
733 e0162726. <https://doi.org/10.1371/journal.pone.0162726>
- 734 Lidén, G. (1969). The scope and application of current audiometric tests. *The Journal of Laryngology*  
735 *& Otology*, 83(6), 507-520.
- 736 Lin, H. W., Furman, A. C., Kujawa, S. G., & Liberman, M. C. (2011). Primary neural degeneration in the  
737 Guinea pig cochlea after reversible noise-induced threshold shift. *Journal of the Association*  
738 *for Research in Otolaryngology : JARO*, 12(5), 605-616. [https://doi.org/10.1007/s10162-011-](https://doi.org/10.1007/s10162-011-0277-0)  
739 [0277-0](https://doi.org/10.1007/s10162-011-0277-0)
- 740 Lobarinas, E., Salvi, R., & Ding, D. (2013). Insensitivity of the audiogram to carboplatin induced inner  
741 hair cell loss in chinchillas. *Hearing Research*, 302, 113-120.  
742 <https://doi.org/10.1016/j.heares.2013.03.012>
- 743 Lopes, A. C., Otubo, K. A., Basso, T. C., Marinelli, É. J. I., & Lauris, J. R. P. (2009). Occupational Hearing  
744 Loss: Tonal Audiometry X High Frequencies Audiometry Perda Auditiva Ocupacional:  
745 Audiometria Tonal X Audiometria de Altas Frequencias.
- 746 Luts, H., Jansen, S., Dreschler, W., & Wouters, J. (2014). Development and normative data for the  
747 Flemish/Dutch Matrix test.

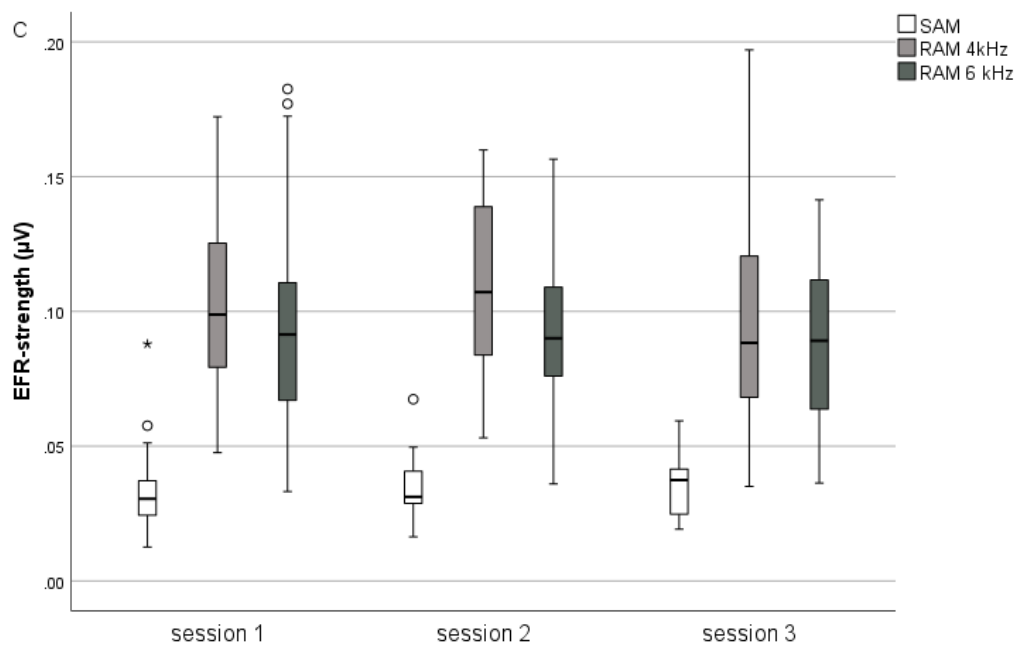
- 748 Maele, T. V., Keshishzadeh, S., Poortere, N. D., Dhooge, I., Keppler, H., & Verhulst, S. (2021). The  
749 variability in potential biomarkers for cochlear synaptopathy after recreational noise  
750 exposure. *Journal of Speech, Language, and Hearing Research*, *64*(12), 4964-4981.
- 751 Mehraei, G., Hickox, A. E., Bharadwaj, H. M., Goldberg, H., Verhulst, S., Charles Liberman, M., &  
752 Shinn-Cunningham, B. G. (2016). Auditory brainstem response latency in noise as a marker  
753 of cochlear synaptopathy [Article]. *Journal of Neuroscience*, *36*(13), 3755-3764.  
754 <https://doi.org/10.1523/JNEUROSCI.4460-15.2016>
- 755 Mehrparvar, A. H., Mirmohammadi, S. J., Ghoreyshi, A., Mollasadeghi, A., & Loukzadeh, Z. (2011).  
756 High-frequency audiometry: a means for early diagnosis of noise-induced hearing loss. *Noise*  
757 *& health*, *13*(55), 402-406. <https://doi.org/10.4103/1463-1741.90295>
- 758 Mepani, A. M., Verhulst, S., Hancock, K. E., Garrett, M., Vasilkov, V., Bennett, K., de Gruttola, V.,  
759 Liberman, M. C., & Maison, S. F. (2021). Envelope following responses predict speech-in-  
760 noise performance in normal-hearing listeners. *Journal of Neurophysiology*, *125*(4), 1213-  
761 1222.
- 762 Mills, D. M., Feeney, M. P., Drake, E. J., Folsom, R. C., Sheppard, L., & Seixas, N. S. (2007). Developing  
763 standards for distortion product otoacoustic emission measurements. *The Journal of the*  
764 *Acoustical Society of America*, *122*(4), 2203-2214.
- 765 Mitchell, C., Phillips, D. S., & Trune, D. R. (1989). Variables affecting the auditory brainstem response:  
766 audiogram, age, gender and head size. *Hearing Research*, *40*(1-2), 75-85.
- 767 Moran, M. D. (2003). Arguments for rejecting the sequential Bonferroni in ecological studies. *Oikos*,  
768 *100*(2), 403-405.
- 769 Munjal, S., Panda, N., & Pathak, A. (2016). Long term test-retest reliability of auditory brainstem  
770 response (ABR) and middle latency response (MLR). *Glob J Oto*, *1*, 555-559.
- 771 Ng, I. H.-Y., & Mcpherson, B. (2005). Test-retest reliability of distortion product otoacoustic  
772 emissions in the 1 to 7 kHz range. *Audiological Medicine*, *3*(2), 108-115.
- 773 Oyler, R. F., Lauter, J., & Matkin, N. (1991). Intrasubject variability in the absolute latency of the  
774 auditory brainstem response. *Journal of the American Academy of Audiology*, *2*(4), 206-213.
- 775 Parthasarathy, A., & Kujawa, S. G. (2018). Synaptopathy in the aging cochlea: Characterizing early-  
776 neural deficits in auditory temporal envelope processing [Article]. *Journal of Neuroscience*,  
777 *38*(32), 7108-7119. <https://doi.org/10.1523/JNEUROSCI.3240-17.2018>
- 778 Plack, C. J., Barker, D., & Prendergast, G. (2014). Perceptual consequences of “hidden” hearing loss.  
779 *Trends in hearing*, *18*, 2331216514550621.
- 780 Plack, C. J., Léger, A., Prendergast, G., Kluk, K., Guest, H., & Munro, K. J. (2016). Toward a diagnostic  
781 test for hidden hearing loss. *Trends in hearing*, *20*, 2331216516657466.
- 782 Popelka, G. R., Osterhammel, P. A., Nielsen, L. H., & Rasmussen, A. N. (1993). Growth of distortion  
783 product otoacoustic emissions with primary-tone level in humans. *Hearing Research*, *71*(1-  
784 2), 12-22.
- 785 Prendergast, G., Tu, W., Guest, H., Millman, R. E., Kluk, K., Couth, S., Munro, K. J., & Plack, C. J.  
786 (2018). Supra-threshold auditory brainstem response amplitudes in humans: Test-retest  
787 reliability, electrode montage and noise exposure [Article]. *Hearing Research*, *364*, 38-47.  
788 <https://doi.org/10.1016/j.heares.2018.04.002>
- 789 Rabinowitz, P., Taiwo, O., Sircar, K., Aliyu, O., & Slade, M. (2006). Physician hearing loss. *American*  
790 *Journal of Otolaryngology*, *27*(1), 18-23.
- 791 Rasetshwane, D. M., Argenyi, M., Neely, S. T., Kopun, J. G., & Gorga, M. P. (2013). Latency of tone-  
792 burst-evoked auditory brain stem responses and otoacoustic emissions: level, frequency,  
793 and rise-time effects. *The Journal of the Acoustical Society of America*, *133*(5), 2803-2817.  
794 <https://doi.org/10.1121/1.4798666>
- 795 Reavis, K. M., McMillan, G. P., Dille, M. F., & Konrad-Martin, D. (2015). Meta-analysis of distortion  
796 product otoacoustic emission retest variability for serial monitoring of cochlear function in  
797 adults. *Ear and Hearing*, *36*(5), e251.

- 798 Rhode, W. S., & Smith, P. H. (1985). Characteristics of tone-pip response patterns in relationship to  
799 spontaneous rate in cat auditory nerve fibers. *Hearing Research*, *18*(2), 159-168.
- 800 Rodríguez Valiente, A., Trinidad, A., García Berrocal, J., Górriz, C., & Ramírez Camacho, R. (2014).  
801 Extended high-frequency (9–20 kHz) audiometry reference thresholds in 645 healthy  
802 subjects. *International Journal of Audiology*, *53*(8), 531-545.
- 803 Schaette, R., & McAlpine, D. (2011). Tinnitus with a normal audiogram: physiological evidence for  
804 hidden hearing loss and computational model. *The Journal of neuroscience : the official*  
805 *journal of the Society for Neuroscience*, *31*(38), 13452-13457.  
806 <https://doi.org/10.1523/JNEUROSCI.2156-11.2011>
- 807 Schlauch, R. S., & Carney, E. (2011). Are false-positive rates leading to an overestimation of noise-  
808 induced hearing loss?
- 809 Schmuziger, N., Patscheke, J., & Probst, R. (2007). An assessment of threshold shifts in  
810 nonprofessional pop/rock musicians using conventional and extended high-frequency  
811 audiometry. *Ear and Hearing*, *28*(5), 643-648.
- 812 Schmuziger, N., Probst, R., & Smurzynski, J. (2004). Test-retest reliability of pure-tone thresholds  
813 from 0.5 to 16 kHz using Sennheiser HDA 200 and Etymotic Research ER-2 earphones. *Ear*  
814 *and Hearing*, *25*(2), 127-132.
- 815 Singh, R., Saxena, R., & Varshney, S. (2009). Early detection of noise induced hearing loss by using  
816 ultra high frequency audiometry. *Int J Otorhinolaryngol*, *10*(2), 1-5.
- 817 Sininger, Y., & Cone-Wesson, B. (2002). Threshold prediction using ABR and SSEPs with infant and  
818 young children. *JackKatz. Handbook of clinical audiology*, 307-321.
- 819 Skoe, E., Camera, S., & Tufts, J. (2019). Noise exposure may diminish the musician advantage for  
820 perceiving speech in noise. *Ear and Hearing*, *40*(4), 782-793.
- 821 Skoe, E., & Tufts, J. (2018). Evidence of noise-induced subclinical hearing loss using auditory  
822 brainstem responses and objective measures of noise exposure in humans [Article]. *Hearing*  
823 *Research*, *361*, 80-91. <https://doi.org/10.1016/j.heares.2018.01.005>
- 824 Swanepoel, D. W., Mngemane, S., Molemong, S., Mkwanazi, H., & Tutshini, S. (2010). Hearing  
825 assessment—reliability, accuracy, and efficiency of automated audiometry. *Telemedicine*  
826 *and e-Health*, *16*(5), 557-563.
- 827 Taberner, A. M., & Liberman, M. C. (2005). Response properties of single auditory nerve fibers in the  
828 mouse. *Journal of Neurophysiology*, *93*(1), 557-569.
- 829 Van Der Biest, H., Keshishzadeh, S., Keppler, H., Dhooge, I., & Verhulst, S. (2023). Envelope following  
830 responses for hearing diagnosis: Robustness and methodological considerations. *The Journal*  
831 *of the Acoustical Society of America*, *153*(1), 191-208.
- 832 Vasilkov, V., Garrett, M., Mauermann, M., & Verhulst, S. (2021). Enhancing the sensitivity of the  
833 envelope-following response for cochlear synaptopathy screening in humans: The role of  
834 stimulus envelope. *Hearing Research*, *400*, 108132.
- 835 Verhulst, S., Ernst, F., Garrett, M., & Vasilkov, V. (2018). Suprathreshold psychoacoustics and  
836 envelope-following response relations: Normal-hearing, synaptopathy and cochlear gain  
837 loss. *Acta Acustica United with Acustica*, *104*(5), 800-803.
- 838 Verhulst, S., Jagadeesh, A., Mauermann, M., & Ernst, F. (2016). Individual Differences in Auditory  
839 Brainstem Response Wave Characteristics: Relations to Different Aspects of Peripheral  
840 Hearing Loss [Article]. *Trends in hearing*, *20*.  
841 [http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L6213725](http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L621372574)  
842 [74](http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L621372574)
- 843 Wagner, W., Heppelmann, G., Vonthein, R., & Zenner, H. P. (2008). Test–retest repeatability of  
844 distortion product otoacoustic emissions. *Ear and Hearing*, *29*(3), 378-391.
- 845 Wang, Y., Yang, B., Li, Y., Hou, L., Hu, Y., & Han, Y. (2000). Application of extended high frequency  
846 audiometry in the early diagnosis of noise--induced hearing loss. *Zhonghua Er Bi Yan Hou Ke*  
847 *Za Zhi*, *35*(1), 26-28.

- 848 Zhao, F., & Stephens, D. (1999). Test-retest variability of distortion-product otoacoustic emissions in  
849 human ears with normal hearing. *Scandinavian Audiology*, 28(3), 171-178.
- 850 Zhu, L., Bharadwaj, H., Xia, J., & Shinn-Cunningham, B. (2013). A comparison of spectral magnitude  
851 and phase-locking value analyses of the frequency-following response to complex tones. *The*  
852 *Journal of the Acoustical Society of America*, 134(1), 384-395.
- 853

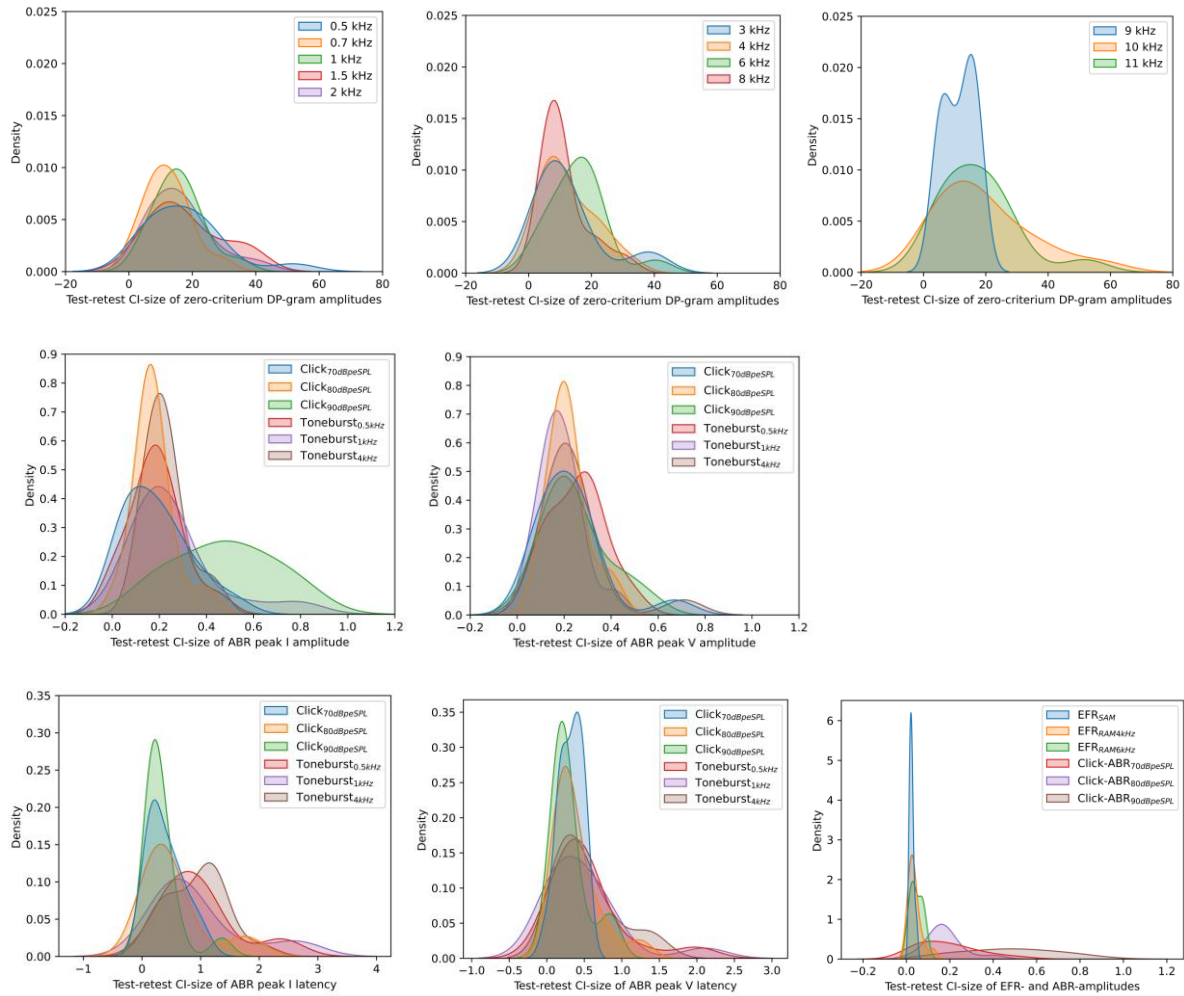
**Figure 1:** Boxplots illustrating speech audiometry results in quiet (A) and noise (B) across sessions. White boxplots indicate SRTs for BB-stimuli, while grey boxplots represent SRTs for HP-conditions. Additionally, panel (C) displays EFR-strengths across sessions, with white, light grey, and dark grey corresponding to SAM-, 4 kHz RAM-, and 6 kHz RAM-EFR strengths, respectively.



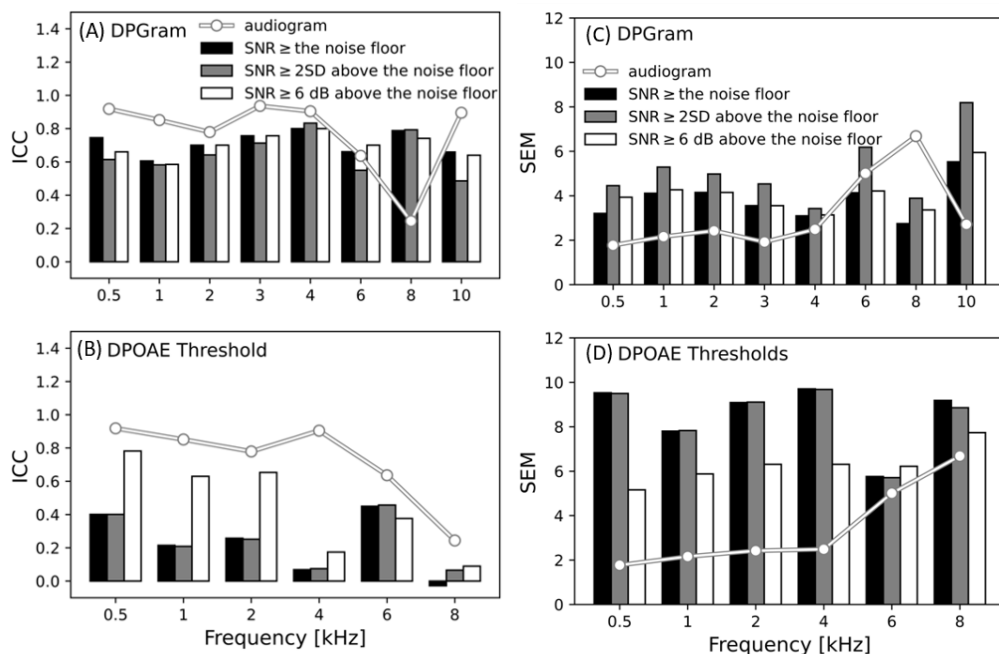




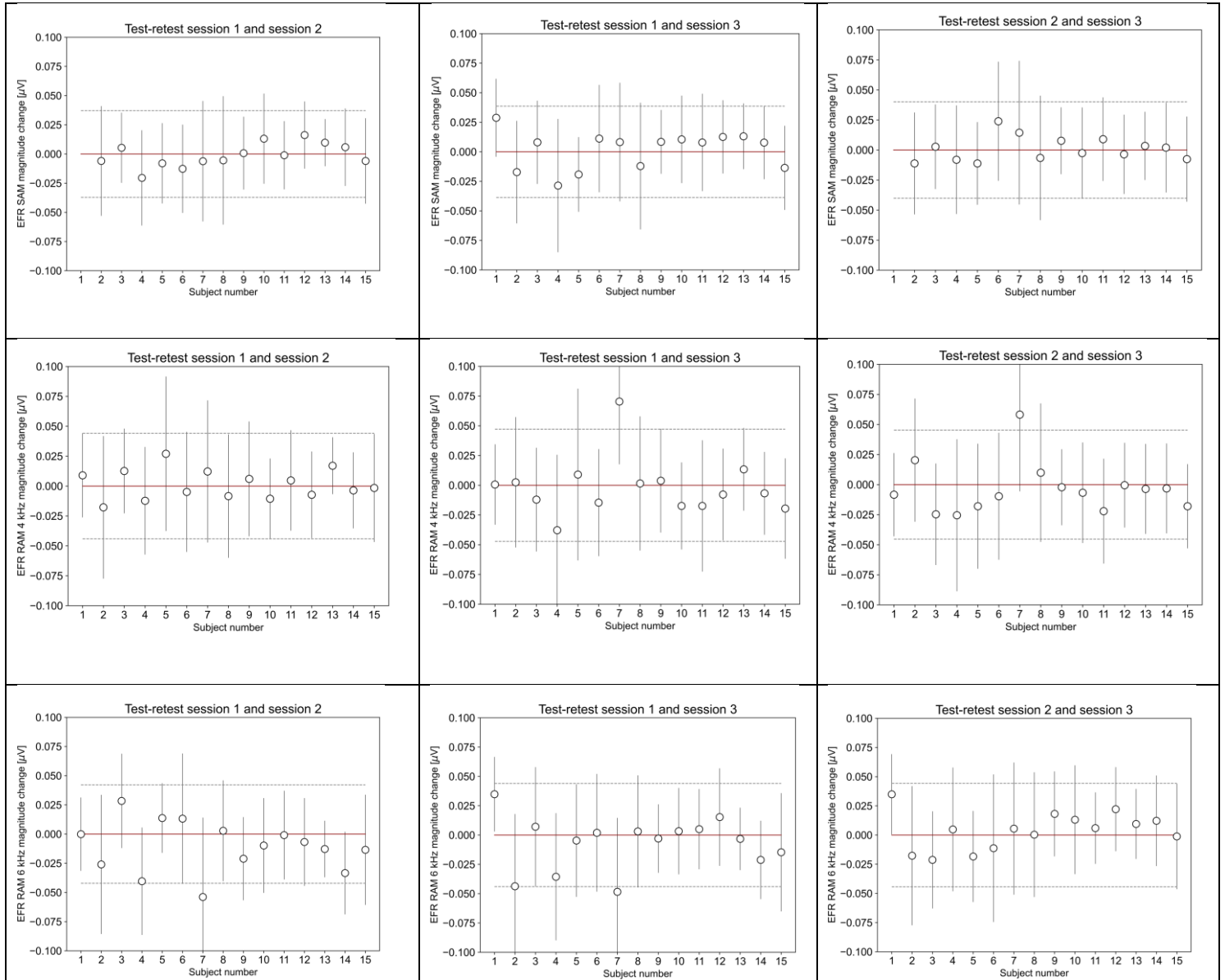
**Figure 2:** Kernel Density Estimate plots illustrating the individual 95% confidence intervals (CIs) for zero-criterion DP-grams across frequency ranges: 0.5-2 kHz (A), 3-8 kHz (B), and 9-11 kHz (C). Panel D-G show ABR amplitudes of peak I (D) and V (E), and ABR latencies of peak I (F) and peak V (G). Additionally, panel (H) displays individual 95% confidence intervals (CIs) for EFR strengths in contrast to click-ABR amplitudes.



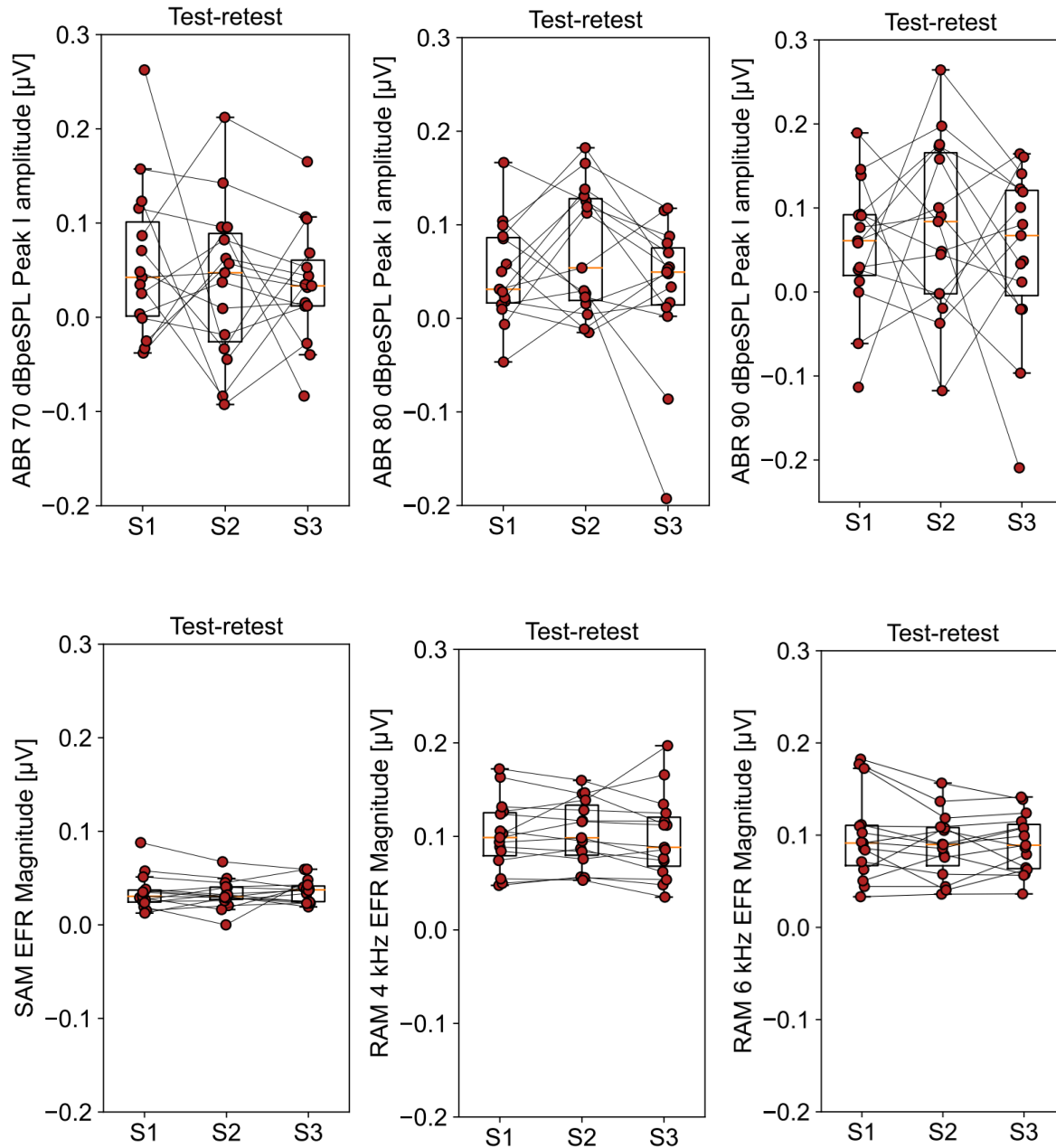
**Figure 3:** Intraclass correlation coefficients (A and B) and standard error of measurements (C and D) of DP-grams and input-output functions in relation to pure-tone audiometry, are illustrated in panels A, B, C and D, respectively. ICCs and SEMs of the three different inclusion criteria regarding signal to noise ratio, i.e.  $SNR \geq$  the noise floor,  $SNR \geq 2SD$  above the noise floor,  $SNR \geq 6$  dB above the noise floor, are illustrated in black, grey and white, respectively, while ICCs and SEMs of pure-tone audiometry are illustrated by a dotted line.



**Figure 4:** Individual EFR-changes between the respective test sessions, represented by 95% confidence intervals (CI) of individual EFR-strength datapoints, and the average CI across all measurements, illustrated by grey dashed lines, for SAM- (row 1), RAM 4 kHz- (row 2) and RAM 6 kHz-stimuli (row 3).



**Figure 5:** Individual ABR-amplitudes and distribution plots for click-ABRs at 70 dBpeSPL, 80 dBpeSPL, and 90 dBpeSPL (row 1), along with individual EFR-magnitudes and distribution boxplots for EFR SAM, RAM 4 kHz, and RAM 6 kHz (row 2). Horizontal lines within the boxplots denote the median ABR-amplitudes and EFR-magnitudes.



S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively.

**Table 1.** Averages per session and frequency, as well as intraclass correlation coefficients (ICC) and standard error of measurements (SEM) for all tested frequencies. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

Frequency (kHz)	0.125	0.250	0.500	1.0	2.0	3.0	4.0	6.0	8.0	10.0	12.5.0	14.0	16.0	20.0
<b>Mean (SD) S1</b>	8.33 (5.876)	4.67 (4.419)	1.67 (5.563)	1.00 (3.873)	3.00 (5.606)	6.43 (8.419)	3.67 (7.432)	13.67 (8.958)	9.00 (6.866)	2.00 (9.024)	4.00 (10.556)	3.33 (13.844)	7.67 (13.478)	5.00 (4.226)
<b>Mean (SD) S2</b>	6.33 (7.188)	4.00 (6.036)	2.00 (7.020)	0.00 (7.071)	2.67 (4.952)	5.00 (8.018)	2.33 (7.528)	14.33 (7.037)	10.33 (8.958)	4.00 (7.838)	3.00 (10.316)	2.33 (13.345)	8.00 (11.148)	3.33 (4.498)
<b>Mean (SD) S3</b>	5.00 (4.629)	0.67 (5.936)	0.67 (6.230)	-0.67 (5.627)	0.00 (4.629)	3.67 (6.673)	1.33 (9.348)	9.67 (8.550)	8.33 (7.480)	2.00 (8.619)	1.00 (9.856)	1.33 (13.157)	7.33 (11.629)	3.67 (7.898)
<b>ICC</b>	0.759***	0.805***	0.918***	0.851***	0.779***	0.936**	0.904***	0.636**	0.244	0.895***	0.911***	0.967***	0.963***	0.087
<b>SEM</b>	2.952	2.504	1.770	2.157	2.418	1.192	2.485	5.008	6.681	2.711	3.010	2.392	2.280	5.447

*S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively*

**Table 2.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of speech audiometry in quiet and noise for BB and HP filtered conditions. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	Speech audiometry in quiet		Speech audiometry in noise	
	BB	HP	BB	HP
<b>Mean (SD) S1</b>	-27.07 (4.011)	29.48 (2.906)	-8.93 (0.806)	-4.01 (1.522)
<b>Mean (SD) S2</b>	-26.31 (4.794)	27.76 (1.763)	-9.15 (0.783)	-4.09 (1.091)
<b>Mean (SD) S3</b>	-28.07 (4.036)	26.64 (2.166)	-9.22 (1.012)	-5.22 (1.424)
<b>ICC</b>	0.341***	0.666***	0.619*	0.694***
<b>SEM</b>	3.470	1.481	0.532	0.797

*S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively*

**Table 3.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of DP-grams. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	500 Hz	700 Hz	1 kHz	1.5 kHz	2 kHz	3 kHz	4 kHz	6 kHz	8 kHz	9 kHz	10 kHz	11 kHz
<b><u>SNR ≥ 0 criterium</u></b>												
<b>Mean (SD) S1</b>	2.85 (5.384)	3.67 (7.238)	6.89 (7.883)	8.89 (6.987)	3.30 (8.157)	-0.29 (8.008)	-7.13 (6.720)	-6.80 (5.809)	-6.68 (6.245)	-10.20 (5.058)	-10.34 (6.556)	-6.47 (6.459)
<b>Mean (SD) S2</b>	1.34 (7.532)	2.89 (6.738)	7.38 (7.609)	9.21 (6.581)	6.03 (8.184)	2.43 (7.731)	-6.03 (6.610)	-5.01 (7.267)	-6.67 (5.387)	-9.163 (5.320)	-6.27 (11.042)	-6.06 (9.883)
<b>Mean (SD) S3</b>	2.01 (6.291)	1.68 (4.524)	7.82 (3.779)	10.46 (6.724)	5.55 (6.498)	3.61 (5.543)	-4.09 (7.522)	-3.93 (8.206)	-8.785 (6.371)	-9.54 (6.816)	-4.95 (10.011)	-3.79 (9.353)
<b>ICC</b>	0.745**	0.774***	0.606*	0.507	0.700**	0.757***	0.800***	0.661**	0.788***	0.762**	0.659**	0.687**
<b>SEM</b>	3.204	2.943	4.111	4.666	4.148	3.552	3.095	4.133	2.746	2.650	5.525	4.800
<b><u>SNR ≥ 2SD criterium</u></b>												
<b>Mean (SD) S1</b>	1.66 (6.142)	-1.16 (9.824)	5.96 (9.194)	8.53 (7.815)	2.68 (9.393)	-1.672 (10.850)	-8.05 (8.026)	-8.41 (7.327)	-8.09 (8.095)	-14.38 (7.687)	-13.77 (8.375)	-10.12 (9.751)
<b>Mean (SD) S2</b>	-2.74 (8.221)	0.75 (8.565)	5.53 (10.348)	9.21 (6.581)	6.03 (8.184)	2.43 (7.731)	-7.39 (8.537)	-6.67 (9.513)	-10.57 (9.196)	-12.93 (7.663)	-8.97 (12.646)	-7.969 (11.359)
<b>Mean (SD) S3</b>	-1.01 (6.755)	-48 (5.524)	7.25 (4.146)	10.46 (6.724)	4.91 (7.494)	3.61 (5.543)	-4.99 (8.830)	-6.65 (10.988)	-12.14 (8.397)	-12.38 (8.502)	-8.65 (12.745)	-6.2740 (11.570)
<b>ICC</b>	0.614*	0.678**	0.582*	0.445	0.642*	0.713**	0.833***	0.550*	0.793***	0.622*	0.486	0.740**
<b>SEM</b>	4.453	4.557	5.291	5.175	4.978	4.530	3.425	6.184	3.888	4.810	8.190	5.502
<b><u>SNR ≥ 6dB criterium</u></b>												
<b>Mean (SD) S1</b>	1.82 (5.524)	3.30 (7.746)	6.77 (8.073)	8.53 (7.815)	3.297 (8.157)	-0.29 (8.008)	-7.13 (6.720)	-7.04 (6.205)	-6.94 (6.574)	-11.26 (6.273)	-10.94 (7.386)	-6.94 (6.977)
<b>Mean (SD) S2</b>	0.17 (8.263)	2.31 (7.389)	7.38 (7.609)	9.21 (6.581)	6.03 (8.184)	2.43 (7.731)	-6.03 (6.610)	-5.30 (7.933)	-6.93 (5.857)	-9.73 (6.244)	-6.83 (11.487)	-7.12 (10.847)
<b>Mean (SD) S3</b>	1.55 (6.508)	0.89 (5.166)	7.82 (3.780)	10.46 (6.724)	5.55 (6.498)	3.61 (5.5.43)	-4.30 (7.835)	-4.67 (9.030)	-9.76 (7.350)	-10.557 (7.150)	-5.08 (10.152)	-4.25 (9.866)
<b>ICC</b>	0.660**	0.811***	0.585*	0.445	0.700**	0.757***	0.800***	0.700**	0.741**	0.788***	0.640**	0.644*
<b>SEM</b>	3.931	2.947	4.268	5.175	4.148	3.552	3.136	4.216	3.363	2.969	5.949	5.522

**Table 4.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of input-output functions. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	500 Hz	1000 Hz	2000 Hz	4000 Hz	6000 Hz	8000 Hz
<b>SNR ≥ 0 criterium</b>						
<b>Mean (SD) S1</b>	12.83 (12.635)	13.69 (9.107)	17.83 (11.308)	29.45 (6.256)	25.740 (6.868)	25.65 (8.190)
<b>Mean (SD) S1</b>	18.14 (11.932)	11.61 (9.739)	16.65 (11.820)	22.84 (12.815)	28.13 (6.762)	21.68 (8.456)
<b>Mean (SD) S1</b>	13.98 (12.893)	14.41 (7.766)	16.074 (8.831)	28.89 (9.571)	27.78 (9.747)	28.16 (9.917)
<b>ICC</b>	0.401	0.215	0.258	0.069	0.450	-0.028
<b>SEM</b>	9.529	7.806	9.089	9.706	5.762	9.186
<b>SNR ≥ 2SD criterium</b>						
<b>Mean (SD) S1</b>	12.83 (12.635)	13.69 (12.635)	19.86 (8.445)	29.45 (6.256)	25.74 (6.868)	25.65 (8.190)
<b>Mean (SD) S1</b>	18.14 (11.932)	11.61 (11.932)	16.65 (11.820)	22.84 (12.815)	28.126 (6.762)	21.68 (8.456)
<b>Mean (SD) S1</b>	13.80 (9.917)	14.36 (7.748)	16.22 (8.785)	28.93 (9.590)	28.06 (9.677)	29.21 (9.841)
<b>ICC</b>	0.401	0.209	0.252	0.075	0.457	0.066
<b>SEM</b>	9.499	7.830	9.113	9.683	5.712	8.856
<b>SNR ≥ 6 dB criterium</b>						
<b>Mean (SD) S1</b>	15.72 (12.208)	14.177 (9.677)	20.50 (8.799)	30.126 (6.824)	27.24 (6.629)	29.66 (5.045)
<b>Mean (SD) S2</b>	18.36 (11.061)	13.81 (11.078)	20.71 (11.442)	29.12 (7.696)	28.40 (7.048)	25.96 (9.345)
<b>Mean (SD) S3</b>	19.57 (9.841)	16.36 (8.517)	17.05 (9.262)	27.05 (9.262)	28.57 (10.111)	29.40 (9.303)
<b>ICC</b>	0.782	0.630	0.653	0.175	0.377	0.090
<b>SEM</b>	5.163	5.877	6.309	6.313	6.223	7.740

*S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively*



**Table 5.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of Click-ABR amplitudes for 70, 80 and 90 dBpeSPL. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	70 dBpeSPL		80 dBpeSPL		90 dBpeSPL	
	Latency	Amplitude	Latency	Amplitude	Latency	Amplitude
<b>Peak I</b>						
Mean (SD) S1	3.37 (0.315)	0.06 (0.082)	3.02 (0.444)	0.05 (0.053)	2.75 (0.210)	0.06 (0.078)
Mean (SD) S2	3.40 (0.251)	0.04 (0.084)	3.09 (0.283)	0.07 (0.068)	2.69 (0.233)	0.08 (0.104)
Mean (SD) S3	3.35 (0.340)	0.04 (0.062)	2.99 (0.316)	0.03 (0.080)	2.76 (0.124)	0.05 (0.103)
ICC	0.640*	0.273	0.787***	0.518*	0.755**	0.100
SEM	0.179	0.064	0.161	0.048	0.096	0.090
<b>Peak III</b>						
Mean (SD) S1	5.63 (0.277)	0.11 (0.109)	5.27 (0.258)	0.06 (0.096)	5.00 (0.227)	0.13 (0.106)
Mean (SD) S2	5.68 (0.353)	0.12 (0.052)	5.25 (0.323)	0.11 (0.096)	4.96 (0.241)	0.12 (0.127)
Mean (SD) S3	5.66 (0.350)	0.07 (0.084)	5.31 (0.333)	0.09 (0.093)	4.89 (0.209)	0.09 (0.109)
ICC	0.967***	0.287	0.968***	0.647**	0.703**	0.552*
SEM	0.058	0.073	0.054	0.057	0.123	0.076
<b>Peak V</b>						
Mean (SD) S1	7.43 (0.394)	0.26 (0.076)	7.17 (0.308)	0.31 (0.117)	6.86 (0.309)	0.32 (0.135)
Mean (SD) S2	7.46 (0.346)	0.24 (0.118)	7.05 (0.325)	0.27 (0.123)	6.93 (0.279)	0.27 (0.137)
Mean (SD) S3	7.47 (0.346)	0.30 (0.118)	7.13 (0.306)	0.32 (0.148)	6.83 (0.264)	0.37 (0.120)
ICC	0.969***	0.691**	0.916***	0.880***	0.937***	0.821***
SEM	0.062	0.059	0.090	0.045	0.070	0.057

S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively

**Table 6.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of TB-ABR amplitudes for 0.5 kHz, 1 kHz and 4 kHz. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	0.5 kHz		1 kHz		4 kHz	
	Latency	Amplitude	Latency	Amplitude	Latency	Amplitude
<b>Peak I</b>						
Mean (SD) S1	3.86 (0.620)	0.08 (0.064)	3.89 (0.597)	0.03 (0.074)	3.31 (0.572)	0.04 (0.069)
Mean (SD) S2	3.76 (0.595)	0.05 (0.098)	3.76 (0.781)	0.04 (0.056)	3.42 (0.502)	0.03 (0.056)
Mean (SD) S3	3.73 (0.564)	0.05 (0.065)	3.71 (0.581)	0.05 (0.083)	3.53 (0.448)	0.01 (0.084)
ICC	0.841***	0.524*	0.859***	0.241	0.824***	0.044
SEM	0.232	0.053	0.244	0.061	0.214	0.072
<b>Peak III</b>						
Mean (SD) S1	6.53 (0.714)	0.06 (0.059)	6.42 (0.406)	0.08 (0.118)	5.91 (0.452)	0.09 (0.082)
Mean (SD) S2	6.35 (0.710)	0.05 (0.0778)	6.21 (0.461)	0.04 (0.081)	5.70 (0.431)	0.09 (0.104)
Mean (SD) S3	6.40 (0.666)	0.07 (0.067)	6.34 (0.397)	0.05 (0.109)	5.80 (0.501)	0.10 (0.103)
ICC	0.843***	0.603*	0.636*	0.331	0.876***	0.687**
SEM	0.272	0.044	0.255	0.086	0.162	0.054
<b>Peak V</b>						
Mean (SD) S1	8.56 (0.451)	0.20 (0.086)	8.16 (0.402)	0.21 (0.081)	7.72 (0.273)	0.23 (0.101)
Mean (SD) S2	8.41 (0.439)	0.21 (0.096)	8.06 (0.401)	0.25 (0.100)	7.55 (0.327)	0.26 (0.100)
Mean (SD) S3	8.53 (0.459)	0.23 (0.072)	8.09 (0.339)	0.25 (0.091)	7.64 (0.180)	0.27 (0.129)
ICC	0.910*	0.601**	0.836***	0.778***	0.729***	0.679**
SEM	0.133	0.053	0.152	0.043	0.140	0.062

S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively

**Table 7.** Intraclass correlation coefficients (ICC) and standard error of measurements (SEM) of EFR-strengths for EFR-SAM and EFR-RAM 4 and 6 kHz stimuli. P-values for between-subjects variability are reflected as \* (0.05<p<0.01), \*\* (0.01<p<0.001), and \*\*\* (p<0.001).

	<b>SAM</b>	<b>RAM 4 kHz</b>	<b>RAM 6 kHz</b>
<b>Mean (SD) S1</b>	0.03 (0.019)	0.10 (0.038)	0.10 (0.047)
<b>Mean (SD) S2</b>	0.03 (0.015)	0.10 (0.035)	0.09 (0.035)
<b>Mean (SD) S3</b>	0.04 (0.012)	0.010 (0.045)	0.09 (0.032)
<b>ICC</b>	0.882***	0.950***	0.930***
<b>SEM</b>	0.005	0.009	0.010

*S1, S2, S3 represent Session 1, Session 2, and Session 3, respectively*