

# Automated cortical auditory response detection strategy

2

3 Fabrice Bardy PhD <sup>a,b</sup>

4 Bram Van Dun PhD <sup>a,c</sup>

5 Mark Seeto <sup>a,c</sup>

6 Harvey Dillon PhD<sup>a,d,e</sup>

7

8

9

10

11 <sup>a</sup> HEARing Co-operative Research Centre, Australia

12 <sup>b</sup> University of Auckland, New Zealand

13 <sup>c</sup> National Acoustic Laboratories, NSW, Australia

14 <sup>d</sup> Macquarie University, NSW, Australia

15 <sup>e</sup> University of Manchester, United Kingdom

16

17

18

19 **Contact: Fabrice Bardy**

20 **[fabrice.bardy@auckland.ac.nz](mailto:fabrice.bardy@auckland.ac.nz)**

21

22 M&HS BUILDING 507 - Bldg 507

23 28 PARK AVE

24 GRAFTON

25 Auckland 1023, New Zealand

26 +64 29 123 06 63

27 **Abstract**

28 **Objective**

29 This study describes a new automated strategy to determine the detection status of an  
30 electrophysiological response.

31 **Design**

32 Response, noise and signal-to-noise ratio of the cortical auditory evoked potential (CAEP) were  
33 characterized. Detection rules were defined: when to start testing, when to conduct subsequent  
34 statistical tests using residual noise as an objective criterion, and when to stop testing.

35 **Study sample**

36 Simulations were run to determine optimal parameters on a large combined CAEP data set collected  
37 in 45 normal-hearing adults and 17 adults with hearing loss.

38 **Results**

39 The proposed strategy to detect CAEPs is fully automated. The first statistical test is conducted when  
40 the residual noise level is equal to or smaller than 5.1  $\mu\text{V}$ . The succeeding Hotelling's  $T^2$  statistical tests  
41 are conducted using pre-defined residual noise levels criteria ranging from 5.1 to 1.2  $\mu\text{V}$ . A rule was  
42 introduced allowing to stop testing before the maximum number of recorded epochs is reached,  
43 depending on a minimum p-value criterion.

44 **Conclusion**

45 The proposed framework can be applied to systems which involves detection of electrophysiological  
46 responses in biological systems containing background noise. The proposed detection algorithm which  
47 optimize sensitivity, specificity, and recording time has the potential to be in clinical setting.

48

49 **Keywords**

- 50 • Electrophysiology
- 51 • Objective response detection
- 52 • Automated algorithm
- 53 • Cortical auditory evoked potentials
- 54 • Residual noise level criteria

## 55 **Introduction**

56 A relevant question when recording electrophysiological responses of any kind is to know whether a  
57 response is present, absent or if the recording is inconclusive. This can be evaluated by a human tester  
58 or by an automated algorithm.

59 Automated algorithms have the potential to be more efficient than human testers by avoiding tester  
60 bias, controlling the false positive rate, and reducing recording time by using sophisticated detection  
61 methods. Objective evoked potential (EP) detection can be achieved using different statistical  
62 techniques, e.g. cross-correlation, detection theory or parametric approaches like the Hotelling's  $T^2$   
63 (Golding et al, 2009; Valdes Sosa et al, 1987; Hyde et al, 1998). Technique selection is aimed at  
64 correctly identifying physiologically present responses (i.e., a high sensitivity) and correctly rejecting  
65 non-physiologically present responses (i.e, a high specificity or a low false-positive rate). Given  
66 different recordings have different response and noise characteristics, the number of required  
67 response averages is recording dependent. The likelihood of detecting a response increases with  
68 increasing response size and decreases with increasing background activity of the recorded signal (i.e.  
69 electrical noise). Multiple responses need to be acquired in order to lower residual noise (RN; the  
70 noise in the averaged response) and reach a signal-to-noise ratio with an acceptable likelihood of  
71 response detection (British Society of Audiology, 2016).

72 When evaluating response presence, there is a need to define time intervals between statistical tests  
73 and criteria to stop testing in the case of response absence. To address the question of when to  
74 statistically evaluate response detection, classic approaches include testing at regular time intervals  
75 or performing one single test at the end of a recording with predefined length. These approaches have  
76 some drawbacks. When applying a statistical test only once after the recording is finished, one  
77 encounters the risk that if the chosen recording length is on the short side, small responses will not  
78 be detected (resulting in a lower than optimal sensitivity). If the predefined recording length is long  
79 however, recording times will be unnecessarily extended when response amplitudes are large.

80 Conversely, testing at regular intervals and repeating this until a response is detected, might be a  
81 better option. In this strategy, potentially time can be saved as a response might be detected after  
82 only a few statistical tests. On the other hand, because the residual noise in the averaged response is  
83 assumed to proportionally decrease with recording length (Elberling & Don, 1984), a fixed interval  
84 between two statistical tests will result in smaller and smaller decreases in residual noise with each  
85 succeeding test. As a result, statistical tests are increasingly likely to be conducted without there being  
86 significant improvements in the signal-to-noise ratio. There is therefore only a small chance that the  
87 response detection is more likely than in the previous statistical test. An additional issue with this  
88 strategy is the large number of tests. Multiple testing increases the probability of false rejection of the  
89 null hypothesis (Type II error or a low specificity). Therefore, the larger the number of tests, the stricter  
90 the correction of the p-value needs to be to keep the false positive rate (FPR) at 5%.

91 To address these drawbacks, another approach will be presented in this paper which takes the residual  
92 noise values into consideration. It will be shown that this approach allows a balance between test  
93 interval lengths and the number of statistical tests. This method is adaptive in the sense that the  
94 interval between two statistical tests will vary depending to the noisiness of the tested subject. The  
95 method allows also to control for the number of statistical tests. Appropriate criteria can be derived  
96 through simulations on a large sample of real-life data sets (Stürzebecher et al, 2005). Finally, deciding  
97 when a response is absent is critical, creating the need for an appropriate stopping criterion. This  
98 criterion is generally determined by a maximum number of epochs or by a sufficiently low residual  
99 noise to allow a likely detection of a predefined response amplitude. A comprehensive overview with  
100 guidelines and suggestions for CAEP testing, detection and absence criteria as used by the British  
101 Society of Audiology, a leading body on CAEP testing, can be found in the British Society of Audiology  
102 Recommended Procedure for CAEP testing (2016). They highlight the relevance of the residual noise  
103 level required for response presence and absence.

104 In this article, we will be describing a fully automated algorithm for response detection and its  
105 optimization procedure. The aim is to have an objective detection of electrophysiological responses  
106 with the highest sensitivity and a controlled specificity in the shortest possible recording time.  
107 Although the proposed techniques in this paper can be applied to any type of electrophysiological  
108 response, the idea is conceptualized through real-life data sets involving the recording of cortical  
109 auditory evoked potentials (CAEPs) for objective hearing threshold estimation in adult subjects. Given  
110 some people cannot provide reliable behavioral feedback due to medical reasons (dementia and  
111 stroke), age (babies and young children), or do not want to because of medico-legal situations  
112 (workers compensation), objective measures offer an alternative to behavioral methods. One of these  
113 objective measures are CAEPs, which are electric responses from the auditory cortex which have been  
114 shown to be a good measure to estimate hearing thresholds in adults (Perl et al, 1953; Beagley &  
115 Kellogg, 1969; Pratt & Sohmer, 1978; Coles & Mason, 1984; Ross et al, 1999; Lightfoot & Kennedy,  
116 2006). Therefore, in more specific terms, we will be describing a fully automated hearing threshold  
117 estimation algorithm. In the first part of the paper, we characterize the response of interest. In the  
118 second part, we define the rules of the algorithm and its optimization procedure using simulations  
119 with several real-life data sets.

120 **Materials and Methods**

121 A large combined CAEP data set recorded on adult subjects was used for the simulations. Optimal  
122 parameters for response detection of CAEPs in adults are derived using the proposed strategy.

123 **Subjects and stimuli**

124 The data used for the simulation were collected during four studies conducted at the National Acoustic  
125 Laboratories. Overall, CAEPs were recorded in 45 normal-hearing adults and 17 adults with hearing  
126 loss in response to short pure- and multi-tone auditory stimuli (either 50 or 70 ms) presented  
127 monaurally via insert earphones (Etymotic Research ER-3A). Sensation levels of the stimuli were 10,  
128 20 and 40 dB SL and the stimulus onset asynchrony (SOA) was randomized uniformly between 1000  
129 and 3000 ms. All stimuli were acoustically calibrated at 70 dB HL according to the ISO standard 389-2  
130 (ISO 1994) in an HA-2 2-cc coupler, incorporating a 1-inch 4144 microphone, a 1-to-1/2 inch DB0375  
131 adaptor, and a 4230 sound level meter (all Brüel & Kjær). Non-response epochs were collected using  
132 portion of EEG signal selected randomly between 1 and 1.3 s after any stimulus onset in case the SOA  
133 was higher than 2s. In total 66442 non-response epochs were collected in the 62 subjects. The false  
134 detection rate for the non-response data was 4.9% when conducting a single Hotelling's T2 with a p-  
135 value criterion of 0.05. This confirmed that the non-response data have similar characteristics of true  
136 non-response data. Table 1 summarizes the relevant details of the four studies. All subjects were in  
137 good general health and reported normal neurological status.

138 **CAEP recordings**

139 The EEG recording equipment used in the four studies was a Neuroscan Synamps2 version 4.3  
140 (Compumedics, Charlotte, NC, USA). The EEG was obtained from 3 gold-plated electrodes placed at Cz  
141 (active), the mastoid contralateral to the ear of stimulation (reference) and the forehead as the  
142 common (ground) channel. Electrode impedance was checked before and after each recording, and  
143 kept under 5 kOhms between active and ground, and between reference and ground. During testing,

144 subjects were seated comfortably in a dimmed, sound attenuated booth. Subjects watched a muted  
145 close-captioned DVD of their choice and were instructed to ignore the stimulus being presented in  
146 their ear.

147 All EEG channels were amplified by a factor of 1210, sampled at 1 kHz, and band-pass filtered online  
148 between 0.1 and 30 Hz. The recording window consisted of a 300 ms pre- and 600 ms post-stimulus  
149 interval (900 ms per epoch). Baseline correction was applied to each individual sweep based on the  
150 average over 100 ms prior to stimulus onset. Epochs exceeding  $\pm 75 \mu\text{V}$  were excluded. Matlab  
151 (MathWorks) and the EEGLAB toolbox (Delorme & Makeig, 2004) were used to process the EEG files.

152

### 153 **Results and simulations:**

#### 154 **PART 1: Characterization of the electrophysiological response, noise and SNR at detection**

155 In order to determine optimal response detection parameters, the response and noise properties of  
156 the signal of interest need to be characterised first. In the case of CAEP, an estimate of the CAEP  
157 amplitude as well as the residual noise (RN) amplitude need to be calculated. This allows  
158 determination of the signal-to-noise ratio (SNR) at detection.

#### 159 **Residual noise (RN) amplitude**

160 The rms amplitude of the RN is estimated based on the epoch-to-epoch variation at each and every  
161 point within the epoch (in a region of interest from 51 to 347 ms after onset). That is, at each point in  
162 the epoch, the variance across epochs is calculated. These values are averaged across all such points  
163 in the epoch, and the square root of that average is taken (Elberling & Don, 1984). This estimation  
164 assumes EEG stationarity. Although this assumption is not completely valid, the accuracy of the RN  
165 prediction is sufficient in a practical sense if noise variance is not changing considerably between  
166 epochs. Figure 1a shows the median RN rms amplitudes and standard deviations (SDs) across  
167 participants after averaging a specific number of epochs for a group of normal-hearing adults (Bardy



168 et al, 2015a). A logarithmic decrease of RN inversely proportional to the square root of the number of  
169 epochs can be observed. The mean rms amplitude per epoch, calculated as the RN amplitude after n  
170 epochs (i.e. 70 in this case) multiplied by the square root of number of epochs, was 12.5  $\mu\text{V}$  (SD 2.65  
171  $\mu\text{V}$ ).

## 172 **CAEP amplitude estimation**

173 To obtain an estimate of the CAEP amplitude, a correction is required by accounting for the RN. This  
174 correction can be applied under the assumption of independence between RN and the true CAEP  
175 (Elberling & Don, 1984). An automated estimate of the CAEP amplitude can be calculated by first  
176 subtracting the RN power from the CAEP power which is calculated as the average waveform power  
177 in a region of interest time interval (i.e. from 51 to 347 ms after onset). Then, the CAEP amplitude  
178 estimate is calculated as the square root of this subtraction.

179 Figure 1b shows the CAEP amplitude distributions at 3 sensation levels: 10, 20 and 40 dB SL for normal-  
180 hearing adults (Bardy et al, 2015a). When combining the distributions obtained at 10, 20 and 40 dB  
181 SL, only 15.5% of CAEP peak amplitudes were larger than 5.1  $\mu\text{V}$ .

182 Knowing now both the CAEP amplitude and RN rms amplitude distributions, the only measure which  
183 still needs to be characterised is the signal-to-noise ratio (SNR) required for a response likely to be  
184 detected. If this specific SNR is known, it is then possible to estimate the maximally allowable RN rms  
185 amplitude at which the first statistical test should occur, still guaranteeing a high likelihood to detect  
186 a CAEP.

187

188 [Insert Fig. 1 here]

189

## 190 **Response detection using Hotelling's $T^2$ & control of the false positive rate**

191 One objective measure for detection of CAEP waveforms is the Hotelling's  $T^2$  statistic, which has been  
192 validated in both adults (Golding et al, 2009) and infants (Carter et al, 2010), and which has been  
193 shown to be at least as accurate as human examiners. Several steps were taken before the Hotelling's  
194  $T^2$  was applied. First, each epoch was divided into 9 bins, with each bin covering a predefined latency  
195 range. The 9 bins covered the range from 51 to 347 ms, with each bin being 33 ms wide. The bin width  
196 and number of bins were chosen based on earlier data (Golding et al, 2009). Second, EEG samples in  
197 these bins were averaged. Hence, each epoch was reduced to a 9-dimensional binned epoch, and the  
198 recorded waveform to a N-by-9 matrix with N the number of collected epochs. Finally, for response  
199 detection, a p-value was obtained from a one-sample Hotelling's  $T^2$  test on the N-by-9 matrix, which  
200 tests the null hypothesis that the true mean vector equals the zero vector (i.e., whether the true  
201 cortical response in every bin is equal to zero).

202 To guarantee a 5% FPR when multiple statistical tests were conducted sequentially, non-response data  
203 were used in simulations to derive a statistical p-value criterion for the Hotelling's  $T^2$  statistic. When  
204 only one statistical test is conducted at the end of a recording, it can be shown for real data using  
205 simulations that a p-value of 0.05 corresponds to a FPR of approximately 5%. This is assuming that the  
206 epochs are independent observations from the same multivariate normal distribution. However, as  
207 different – and more complex – algorithms are employed here, simulations need to be conducted to  
208 control the FPR as it is difficult to mathematically derive which p-value needs to be applied.

### 209 **Signal-to-noise ratio at detection**

210 Signal detection depends primarily on the characteristics of signal and noise, both reflected in the SNR  
211 measure. In this section, the SNR is defined first. Then, the SNR at which a response is likely to be  
212 detected using the Hotelling's  $T^2$  statistic is investigated. These characteristics will allow the derivation  
213 of criteria guiding when to conduct statistical tests.

214 The signal-to-noise ratio (in dB) of the CAEP at detection is defined as:

215 
$$\text{SNR (dB)} = 20 \log_{10} \frac{\text{CAEP amplitude}}{\text{RN amplitude}}, \text{ with} \quad (\text{Eq. 1})$$

- 216 - SNR: Signal to noise ratio of the CAEP amplitude and the RN amplitudes;  
217 - CAEP amplitude, as defined in section “CAEP amplitude estimation ”; and  
218 - RN amplitude, based on the epoch-to-epoch standard deviation (see section Residual noise  
219 (RN) amplitude).

220 We determined the median SNR needed for a CAEP to be detected using data reported in Bardy et al  
221 (2015a). Using a sequential test strategy, the p-value was calculated after the collection of nine  
222 epochs and subsequently, every additional two epochs. For response detection, a correction for  
223 multiple testing of the p-value (to 0.006) was derived using non-response data to keep the FPR at  
224 5%. For every test condition in each subject, the SNR at CAEP detection was collected. Figure 2  
225 shows the distribution of SNRs at which a significant CAEP could just be detected. As can be derived  
226 from Figure 2, 50% of CAEPs needed a SNR of 3.3 dB or greater to be detected and 18% of detections  
227 occurred at negative SNRs. While the relationship between SNR and Hotelling’s  $T^2$  is highly  
228 correlated there is variation which depends on the shape of the response and characteristics of the  
229 noise. In fact, for a particular shape of response, and a particular distribution of noise rms values  
230 (calculated across epochs) along the epoch, if the noise rms value was lower at every point along the  
231 epoch by the same proportion, then both Hotelling’s  $T^2$  and SNR would increase. Thus for this type of  
232 variation, there would be a perfectly monotonic relationship between the two. However, Hotelling’s  
233  $T^2$  gives the greatest weight to the time windows that have the best combination of signal amplitude  
234 to noise rms. SNR however, weights all time points the same; only the total signal rms and the total  
235 noise rms matter. So, it’s possible that a change in the signal, or the noise has a different effect on  
236 Hotelling’s  $T^2$  than it does on the SNR measure. For example an increase of the noise rms in a time  
237 window where the signal is zero will have virtually no effect on Hotelling’s  $T^2$ , but will decrease  
238 SNR. A noise increase in a time window where the signal is at a maximum will have cause a large  
239 decrease in Hotelling’s  $T^2$ , but only a small decrease in SNR. So, it’s easy to see that although

240 detection with Hotelling's  $T^2$  generally gets easier as SNR improves (over a wide range of possible  
241 SNRs), there will be variations from this relationship. Consequently, it will sometimes be possible to  
242 detect signals below some criterion SNR (such as 0 dB) while sometimes not being able to detect  
243 them for SNRs above this criterion SNR.

244

245 [Insert Fig. 2 here]

246

## 247 **PART 2: Detection rules within a single stimulus condition**

248 The number of statistical tests conducted during the recording of a single stimulus condition needs to  
249 be limited in order to avoid either an unnecessary increase of the FPR or an excessive decrease of the  
250 p criterion used for each test. The purpose of this section is to define the strategy and rules to be used  
251 in real-time during data collection for response detection. The rules described are based on the  
252 characteristics of the response amplitude, the RN and SNR at detection described in Part 1. The aim  
253 is: 1) to determine the criteria to start statistical testing, 2) to define when to perform successive  
254 statistical tests, and 3) to define when to stop collecting data. Finally, the validity of the method is  
255 demonstrated using real data through simulations.

### 256 **When to start statistical testing for response detection?**

257 First, the minimum number of epochs to conduct the first statistical test needs to be larger than the  
258 number of bins to calculate the Hotelling's  $T^2$ . Second, the RN level for the first statistical test is data  
259 driven and depends on the RN at which there already is a reasonable chance to detect a CAEP. From  
260 the data displayed in Figure 1b, it was calculated that 86% of true CAEP peak amplitudes of the  
261 response detected tend to be smaller than 5.1  $\mu$ V. Hence, testing at RN higher than 5.1  $\mu$ V would only  
262 allow appropriate response detection conditions for a minority of CAEPs. This results in the waste of

263 one (or several) statistical tests at the early stages of the recording. Considering these data, one  
264 criterion for the first statistical test to be conducted is a RN amplitude below 5.1  $\mu$ V.

265 **When to conduct statistical tests later on? Residual noise as an objective criterion.**

266 When the first statistical test has been conducted, the question is when to conduct the remaining  
267 statistical tests. There are two approaches which are commonly used when performing statistical tests  
268 on a recording consisting of multiple epochs:

- 269 1) Apply a statistical test only once, after the collection of a fixed number of epochs; and
- 270 2) Multiple tests at fixed intervals (equidistant epochs).

271 We propose a novel approach that relies on multiple tests at predefined RN amplitudes. This approach  
272 guarantees that the SNR improves by a predefined ratio from the time when the previous statistical  
273 test was conducted (assuming the CAEP is constant in amplitude), allowing an increased chance of  
274 detection. Moreover, it implicitly adapts to the noise condition within each recording (which is highly  
275 different depending on the population tested). For example, in cases of increased noise during the  
276 recording, statistical testing will be automatically postponed until it reaches the predefined RN  
277 criterion. In addition, the number of statistical tests and the space between two statistical tests can  
278 be controlled. Figure 3a represents the strategy that has been derived based on the following  
279 constraints:

280 ***The number of statistical tests and their spacing*** is a trade-off between test duration, detection  
281 sensitivity and clinical applicability. Clinical applicability can be defined as having short test durations  
282 tolerable to the patient and a sufficient number of statistical updates for the clinician. A low number  
283 of statistical tests results in higher (less strict) p-value criteria (typically with  $p = 0.05$  in the extreme  
284 case when there is only 1 statistical test at the end of the recording). However, because of the low  
285 number of statistical tests, test duration will be longer as fewer opportunities are available to end a  
286 recording early (with the maximum test duration achieved in the extreme case of only one statistical

287 test at the end of the recording). Conversely, a higher number of statistical tests results in lower  
288 (stricter) p-value thresholds. Test durations will likely be shorter as more opportunities arise to stop  
289 testing earlier. However, the very strict p-values that are the consequence of a large number of tests  
290 may also delay detection, or in some cases prevent it from occurring. To summarize, when test  
291 duration and detection sensitivity are traded off, an optimal number of statistical tests can be derived,  
292 providing clinical applicability is still acceptable.

293 The statistical tests that are executed need to be distributed over a range. A balance needs to be found  
294 between: (1) keeping the spacing between two neighboring tests as small as possible to allow early  
295 response detection, and (2) keeping the spacing as large as possible such that a test would be  
296 worthwhile, given the adverse effects that additional tests have on either the FDR or the p-value  
297 criterion, or both. The RN needs to have dropped by a significant fraction since the previous test  
298 making it more likely that a response is detected were there to be one present. The following strategy  
299 achieves this balance.

300 Figure 3 shows the proposed spacing of tests in three inter-related ways. Panel (a) shows the residual  
301 noise values at which each successive test is carried out. The first test is carried out when RN equals  
302  $5.1 \mu\text{V}$ . Succeeding tests occur as RN decreases to the values determined by the exponential function  
303 shown. This curve applies irrespective of the actual noisiness of any individual person being tested.  
304 For a typical person with an rms noise level (per epoch) of  $12.5 \mu\text{V}$ , panel (b) shows the total number  
305 of epochs that will have elapsed when each test has been carried out, and panel (c) shows the  
306 corresponding number of epochs between immediately adjacent tests. As expected, the number of  
307 epochs between adjacent tests increases with increasing test number. Note, however, that the  
308 number of epochs between adjacent tests never exceeds 40, which was one of the design goals that  
309 helped determine the function shown in panel (a).

310 Figure 3a shows the RN amplitudes at which to conduct the sequential statistical tests. It is  
311 independent from the subject being tested: some subjects might be less noisy, therefore needing a  
312 lower number of epochs to reach the RN criteria.

313 Figure 3b presents the number of epochs needed to reach the RN amplitudes shown in Figure 3a,  
314 when the RN per epoch is equal to  $12.5 \mu\text{V}$ , which corresponds to the mean RN per epoch in normal  
315 hearing adults. The exponential curves indicate that increasingly more epochs need to be collected  
316 after each sequential statistical test to reach the next RN criterion.

317 Figure 3c shows the differential number of epochs between 2 tests before reaching the next RN  
318 criterion and is based on Figure b. It is clear that with increasing statistical test index, the spacing  
319 distance (expressed in number of epochs) between tests increases as well. To keep this spacing under  
320 control, it was opted to limit the distance to about 40 epochs for practical reasons. This constraint in  
321 turn controlled the maximum slope of the RN graph and the minimum number of statistical tests in  
322 Figure 3a.

323 To recapitulate, the exponential function displayed in Figure 3a dictates when to perform each  
324 statistical test. It is determined by its slope, the number of statistical tests and the minimum RN  
325 amplitude that needs to be reached before the first test can be conducted.

326

327 [Insert Fig. 3 here]

328

### 329 **When to stop averaging?**

330 Stopping criteria are determined by CAEP amplitude distributions at various stimulus levels, RN  
331 estimates, and the objective detection algorithm with FPR. Five stopping criteria can be identified:

332

- 333 1. ***When a response is detected*** - using an objective detection technique like e.g. the Hotelling's  
334  $T^2$  with a predefined detection criterion that has been determined a priori using non-response  
335 data.
- 336 2. ***When the maximum number of epochs in a recording has been acquired.*** The maximum  
337 number of epochs is a trade-off between the maximum acceptable test duration and the  
338 required sensitivity for response detection at low sensation levels. Sensitivity in turn depends  
339 on the RN level of the averaged waveform when the maximum number of epochs has been  
340 collected.
- 341 3. ***When the subject appears too noisy*** to continue the process. This can be identified at any  
342 stage in the recording.
- 343 4. ***When it is clear the required objective criterion will not be reached.*** If the p-value is still above  
344 a specific value after a certain number of epochs, the recording can be stopped immediately.  
345 Simulations will be conducted to determine these p-values.
- 346 5. ***When reaching a predefined minimum RN criterion and no response has been detected.*** We  
347 advise against the use of this criterion as it was noticed during live recordings that this  
348 approach can lead to inappropriate results in some people with genuinely low RN (and CAEP)  
349 amplitudes. A statistical failure to detect a response may be due to insufficient averaging to  
350 achieve the required SNR, rather than the absence of a response.

351 The concept of adaptive stopping criteria has been introduced in previous research (Kelley et al, 2018;  
352 Botella et al, 2006). Moreover, diagnostic evoked potential devices, including most commercially  
353 available auditory brainstem response (ABR) and auditory steady-state response (ASSR) instruments,  
354 utilise automatic stopping rules. The theoretical literature describes several approaches to stop testing  
355 earlier than anticipated. For example, in sequential estimation the sample size to use is not specified  
356 at the start, and instead outcomes are employed to evaluate a predefined stopping rule if sampling  
357 should continue or stop (Kelley et al, 2018; Botella et al, 2006) . In addition, Bayesian statistics have  
358 been used to allow premature stopping of behavioural experiments or clinical trials while keeping the



359 false positive rate constant (Psioda et al, 2018; Komaki & Biswas, 2018; Alcalá-Quintana & García-  
360 PÉREZ, 2005)

### 361 **PART 3: Simulations and validation**

362 According to the rules defined in Part 2, the first simulations were conducted to adjust the statistical  
363 detection criterion (p-value) for multiple testing. Second, parameters were defined for p-values that  
364 do not reach a certain minimum value after a number of collected epochs (allowing early stopping).

#### 365 Simulation 1: p-value correction for multiple testing to control the FPR

366 The general aim of Simulation 1 is to find the p-value which keeps the FPR at 5%. A strict and  
367 conservative estimate can be derived using the Bonferroni-correction, which divides the p-value by  
368 the number of statistical tests. A better approximation can be obtained through Monte Carlo  
369 simulations on EEG data containing no CAEPs, calculating the FPR for a range of p-values while  
370 adhering to the rules defined in Part 2. Simulations were conducted using EEG data collected during 4  
371 different studies described in the Materials and Methods section.

372

373 [Insert Table 1 here]

374

375 A total number of 66442 epochs formed 552 simulated recordings of 120 epochs each. The following  
376 procedure was followed for the Monte Carlo simulation to calculate the p-value criterion needed to  
377 achieve a FPR of 0.05.

- 378 - **for** p varied from 0.0001 to 0.05 in steps of 0.0001
  - 379 ○ **for** each simulated recording out of 552
    - 380 ▪ **for** each epoch ranging from 20 to 120
      - 381 ● **add** the epoch to the grand average
      - 382 ● **if** a predetermined RN amplitude is reached

- 383 ○ conduct a statistical test, providing a p-value  $P$
- 384 ○ if  $P \leq p$ , stop and  $FP = FP + 1$
- 385 ○ **calculate**  $FPR = FP / 552$ .
- 386 ○ if  $FPR < 0.05$ , stop

387 The p-value criterion adopted was the highest p-value tested with  $FPR < 0.05$ . Each simulated  
 388 recording allowed between 3-9 statistical tests (mean = 6.2) (depending on the RN amplitudes that  
 389 have been reached), while satisfying an FPR of 5%. For different maximum numbers of epochs, the  
 390 following p-value criteria were determined (in brackets): 120 (0.0077), 110 (0.0119), 100 (0.0129).

391

392 Rule to stop testing if the p-value is still above a certain value after n epochs

393 When we perform a sequence of statistical tests with each test using all epochs in the run up to that  
 394 time, the sets of epochs used for different tests overlap, so the tests are not independent. Because of  
 395 this non-independence, knowledge of the p-value for a particular test allows us to be sure that after  
 396 a specified number of additional epochs, the new p-value will be in a certain range. It follows that if  
 397 we have set a maximum number of epochs to be recorded, then it will sometimes be possible to know  
 398 before reaching the maximum number of epochs that the p-values from later tests will not be less  
 399 than the critical p-value for detection. In other words, it is sometimes possible to know in advance  
 400 that subsequent tests in a run will not detect a response, in which case the run can be stopped early  
 401 to save time.

402 To state the condition for early stopping we require some notation. Consider a test after  $n$  epochs and  
 403 a later test after  $N$  epochs (so  $N > n$ ), and denote the respective p-values by  $p_n$  and  $p_N$ . Let  $k$  be the  
 404 number of bins, so in our framework we have  $k = 9$ . Let  $\Psi_{\nu_1, \nu_2}$  be the cumulative distribution function  
 405 of an F random variable with degrees of freedom  $\nu_1$  and  $\nu_2$ .

406 The general result is that if  $0 < q < 1$  and

407 
$$p_n > 1 - \Psi_{k,n-k} \left[ \left( \frac{n-k}{N} \right) \left( \left( \frac{n}{N-k} \right) \Psi_{k,N-k}^{-1}(1-q) - \frac{N-n}{k} \right) \right],$$

408 then  $p_N > q$ . An outline of the proof of this result is given in the appendix.

409 If we take  $N$  to be the maximum number of epochs and  $q$  to be the p-value cutoff for detection then  
410  $p_N > q$  means that even after the maximum number of epochs, a response is not detected, so the  
411 testing can be stopped after epoch  $n$  instead of waiting until epoch  $N$ .

412 Table 2 shows the critical p-values for early stopping, assuming a maximum of 120 epochs and a p-  
413 value detection criterion of  $p < 0.01$ , obtained from the expression above with  $k = 9$ ,  $q = 0.01$ ,  $N =$   
414 120 and  $n$  from 102 to 119. The interpretation is, for example, that if the p-value after 110 epochs is  
415 greater than 0.298 then the p-value after 120 epochs cannot be less than 0.01, so the testing can be  
416 stopped after 110 epochs.

417

418 [Insert Table 2 here]

419 **Discussion**

420 This paper presented a general framework to optimize the detection of time-locked evoked potentials.  
421 The optimization allows the determination of when to conduct each statistical test and how long to  
422 test for. It aimed to optimize the inevitable tradeoffs between detection sensitivity, FPR (i.e., the  
423 specificity) and recording time. This framework can be applied to detection of auditory evoked  
424 potentials for threshold estimation.

425 **Comparisons with other automated objective response detection paradigms**

426 While guidelines have been developed with recommended (clinical) criteria for response presence and  
427 absence for CAEP testing (British Society of Audiology, 2016), only a handful studies have described  
428 objective statistical techniques to detect CAEPs when determining hearing thresholds, with an  
429 overview provided by Van Dun et al (2015). They do not provide guidelines however on when to  
430 automatically proceed to the next stimulus level, and on how to maximize sensitivity and minimize  
431 data collection time. The only study that we are aware of that provides elements towards an  
432 automatic objective threshold searching paradigm, is (Elberling & Don, 1984) for auditory brainstem  
433 responses (ABRs). They presented practical guidelines related to ABR and residual noise amplitudes,  
434 statistical measures for detection (using the Fsp), sensitivities and false positive rates. Based on these  
435 parameters, they derived when (or when not) to stop data collection and proceed to the next stimulus  
436 level. Overall, the paradigm presented in this study is the first one of its kind, and it allows automation  
437 of cortical threshold searching. A practical strategy is provided below.

438 **Practical strategy for an automated response detection of CAEPs in adult populations**

- 439
- The first statistical test is conducted when the RN level is equal to or smaller than 5.1  $\mu$ V.
  - Succeeding statistical tests are conducted when specific RN levels are reached, as defined in  
441 Figure 3a.

442 • The p-value detection criterion adopted throughout the strategy is equal to  $p = 0.01$ . This  
443 criterion guarantees a FPR of about 5%.

444 Averaging is stopped when:

- 445 • a response is detected using Hotelling's  $T^2$  with a detection criterion of  $p = 0.01$ ; or
- 446 • the maximum number of 120 epochs has been collected; or
- 447 • a RN level of  $5.1 \mu\text{V}$  will likely not be reached before the maximum number of epochs has  
448 been collected; or
- 449 • the p-value at a specific number of accepted epochs ( $\geq 102$ ) is higher than the critical p-value  
450 presented in Table 2.

#### 451 **Possible limitations of the method**

452 A possible limitation of this study arise from the specific functions shown in Figure 3. Although these  
453 are reasonable choices, we are not aware of any technique that could be used to derive them such  
454 that the combination of sensitivity, specificity and recording time could be mathematically optimized.  
455 A better set of functions might therefore exist.

456 A second limitation is that for the technique to be applied to other populations (such as infants), the  
457 true CAEP peak and RN amplitude distributions and resulting SNRs (Figures 1 and 2) need to be  
458 determined in those populations. Nevertheless, we believe that the general approach will be valid.  
459 However, the maximum number of epochs might need to be adjusted for each population, impacting  
460 the spacing of the statistical tests.

461 A third limitation is the absence of an inconclusive result in case the residual noise levels are too  
462 high after the maximum number of epochs has been reached. Currently the system will indicate that  
463 the CAEP is absent. A new criterion could be added for this specific case. The test could be  
464 categorised as inconclusive instead of absent if the subject's mean noise-per-epoch value is higher  
465 than an established age-appropriate threshold when the maximum number of epochs is reached.  
466 This option might need to be implemented in the algorithm in the future.

467

468 **An optimal detection method for CAEPs in a clinical environment**

469 The proposed approach will make CAEP testing more accessible for clinicians, who generally have to  
470 rely on their own judgement how and when to interpret the cortical waveforms, and when to stop  
471 collecting data. This decision process takes time, which is in short supply in a clinical environment.  
472 Moreover, as clinicians all have their own approaches, it is next to impossible to derive a false positive  
473 rate for each clinician. The general technique described is intended for response detection in any age  
474 group, more specifically those who cannot provide reliable feedback like e.g. infants, young children,  
475 malingers, those with multiple disabilities, who have suffered a stroke or are diagnosed with  
476 dementia. The real-time implementation of the proposed algorithm for clinical use has completed on  
477 electrophysiological hardware system called 'HEARLab' developed at the National Acoustic  
478 Laboratories.

479 **Conclusion**

480 This paper determined how often and when to conduct a statistical test and how long to test for the  
481 detection of time-locked evoked responses. When sufficient data are available to run the required  
482 simulations for determining specific parameters, the proposed framework can be applied to any  
483 system which involves detection of time-locked electrophysiological responses in biological systems  
484 containing background noise. Applications of this approach can be found in auditory, visual or motor  
485 threshold estimation, or basically any automated system which can converge to a predetermined  
486 criterion.

487 **Acknowledgements:**

488 We thank Professor Robert Cowan and the two reviewers for helpful comments on the manuscript.  
489 The experimental data used for this research were collected while the authors worked at the  
490 National Acoustic Laboratories. The authors acknowledge the financial support of the Australian  
491 Department of Health, the NSW Health and the HEARing CRC, established under the Australian

492 Government's Cooperative Research Centres (CRC) Program. The CRC Program supports industry-led  
493 collaborations between industry, researchers and the community. Harvey Dillon acknowledges the  
494 support of the NIHR Manchester Biomedical Research Centre and Macquarie University.

495

496 **Declaration of interest statement**

497 The authors have been involved with the development of HEARLab and its modules, but do not  
498 receive financial benefits from its sale.

499

500

501

502 **References**

- 503 Alcalá-Quintana, R., García-PÉREZ, M.A., 2005. Stopping rules in Bayesian adaptive threshold  
504 estimation. *Spat. Vis.*
- 505 Bardy, F., Van Dun, B., Dillon, H., 2015a. Bigger Is Better: Increasing Cortical Auditory Response  
506 Amplitude Via Stimulus Spectral Complexity. *Ear Hear*, 36(6), p.677–687. Available at:  
507 <http://www.ncbi.nlm.nih.gov/pubmed/26039014>.
- 508 Bardy, F., Sjahalam-King, J., Van Dun, B., Dillon, H., 2015b. Cortical auditory evoked potentials  
509 (CAEPs) in response to multi-tone (MT) stimuli in hearing-impaired adults. *Am. J. Audiol.*, (27),  
510 p.406–415.
- 511 Beagley, H.A., Kellogg, S.E., 1969. A comparison of evoked response and subjective auditory  
512 thresholds. *Int. J. Audiol.*, 8, p.345–353.
- 513 Botella, J., Ximénez, C., Revuelta, J., Suero, M., 2006. Optimization of sample size in controlled  
514 experiments: The CLAST rule. *Behav. Res. Methods*.
- 515 British Society of Audiology, 2016. *Recommended Procedure for CAEP testing*, Available at:  
516 <https://www.thebsa.org.uk/wp-content/uploads/2016/05/Cortical-ERA.pdf>.
- 517 Carter, L., Golding, M., Dillon, H., Seymour, J., 2010. The detection of infant cortical auditory evoked  
518 potentials (CAEPs) using statistical and visual detection techniques. *J. Am. Acad. Audiol.*, 21(5),  
519 p.347–356.
- 520 Coles, R.R.A., Mason, S.M., 1984. The results of cortical electric response audiometry in medico-legal  
521 investigations. *Br. J. Audiol.*, 18(2), p.71–78.
- 522 Delorme, A., Makeig, S., 2004. EEGLAB: An open source toolbox for analysis of single-trial EEG  
523 dynamics including independent component analysis. *J. Neurosci. Methods*, 134(1), p.9–21.
- 524 Van Dun, B., Dillon, H., Seeto, M., 2015. Estimating hearing thresholds in hearing-impaired adults



525 through automatic detection of cortical auditory evoked potentials. *J. Am. Acad. Audiol.*, 26(4),  
526 p.1–14.

527 Elberling, C., Don, M., 1984. Quality estimation of averaged auditory brainstem responses. *Scand.*  
528 *Audiol.*, 13(3), p.187–197.

529 Golding, M., Dillon, H., Seymour, J., Carter, L., 2009. The detection of adult cortical auditory evoked  
530 potentials (CAEPs) using an automated statistic and visual detection. *Int. J. Audiol.*, 48(12),  
531 p.833–842.

532 Hyde, M., Sininger, Y.D., Don, M., 1998. Objective detection and analysis of auditory brainstem  
533 response: an historical perspective. *Semin. Hear.*, 19(1), p.97–117.

534 Kelley, K., Darku, F.B., Chattopadhyay, B., 2018. Accuracy in parameter estimation for a general class  
535 of effect sizes: A sequential approach. *Psychol. Methods*.

536 Komaki, F., Biswas, A., 2018. Bayesian optimal response-adaptive design for binary responses using  
537 stopping rule. *Stat. Methods Med. Res.*

538 Lightfoot, G., Kennedy, V., 2006. Cortical electric response audiometry hearing threshold estimation:  
539 accuracy, speed, and the effects of stimulus presentation features. *Ear Hear.*, 27(5), p.443–456.

540 Perl, E.R., Galambos, R., Gloor, A., 1953. The estimation of hearing threshold by  
541 electroencephalography. *Electroencephalogr. Clin. Neurophysiol.*, 5(4), p.501–512.

542 Pratt, H., Sohmer, H., 1978. Comparison of hearing threshold determined by auditory pathway  
543 electric responses and by behavioural responses. *Int. J. Audiol.*, 17(4), p.285–292.

544 Psioda, M.A., Soukup, M., Ibrahim, J.G., 2018. A practical Bayesian adaptive design incorporating  
545 data from historical controls. *Stat. Med.*

546 Ross, B., Lutkenhoner, B., Pantev, C., Hoke, M., 1999. Frequency-specific threshold determination  
547 with the CERAgram method: basic principle and retrospective evaluation of data. *Audiol. {&}*

548 *Neuro-Otology*, 4(1), p.12–27. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/9873149>.

549 Stürzebecher, E., Cebulla, M., Elberling, C., 2005. Automated auditory response detection: Statistical  
550 problems with repeated testing Evaluación repetida en la detección de respuestas auditivas.  
551 *Int. J. Audiol.*, 44(2), p.110–117. Available at:  
552 <http://www.tandfonline.com/doi/full/10.1080/14992020400029228>.

553 Valdes Sosa, M.J., Bobes, M.A., Perez Abalo, M.C., Perera, M., Carballo, J.A., et al, 1987. Comparison  
554 of auditory-evoked potential detection methods using signal detection theory. *Audiology*,  
555 26(3), p.166–178. Available at: [file:///C:/Bibliografia sobre  
556 Neurociencias/pubneuropdf003/comparison.pdf](file:///C:/Bibliografia%20sobre%20Neurociencias/pubneuropdf003/comparison.pdf).

557

558 **Tables**

559 Table 1: Summary of the EEG data sets used for Simulation 1. NH: normal-hearers, HI: hearing-  
560 impaired.

| Data source          | Adult population | N         | Average number of non-response epochs per subject | Total number of non-response epochs |
|----------------------|------------------|-----------|---|-------------------------------------|
| Bardy et al. 2015a   | NH               | 15        | 928   | 15776                               |
| Bardy et al. (2015b) | HI               | 17        | 917   | 15589                               |
| NAL data set 1       | NH               | 17        | 845   | 14366                               |
| NAL data set 2       | NH               | 13        | 1218  | 20711                               |
|                      | <b>Total</b>     | <b>62</b> | <b>3908</b>                                       | <b>66442</b>                        |

561

562 Table 2: Critical p-values for early stopping, assuming a maximum of 120 epochs and a p-value  
563 detection criterion of  $p < 0.01$ . Testing can be stopped after the number of epochs shown in the table  
564 if the p-value is greater than the corresponding p-value in the table.

|                |            |            |            |            |            |            |            |            |            |
|----------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|
| <b>Epoch</b>   | <b>102</b> | <b>103</b> | <b>104</b> | <b>105</b> | <b>106</b> | <b>107</b> | <b>108</b> | <b>109</b> | <b>110</b> |
| <b>p-value</b> | 0.979      | 0.938      | 0.872      | 0.784      | 0.683      | 0.578      | 0.475      | 0.381      | 0.298      |
| <b>Epoch</b>   | <b>111</b> | <b>112</b> | <b>113</b> | <b>114</b> | <b>115</b> | <b>116</b> | <b>117</b> | <b>118</b> | <b>119</b> |
| <b>p-value</b> | 0.229      | 0.172      | 0.127      | 0.092      | 0.066      | 0.047      | 0.033      | 0.023      | 0.015      |

565

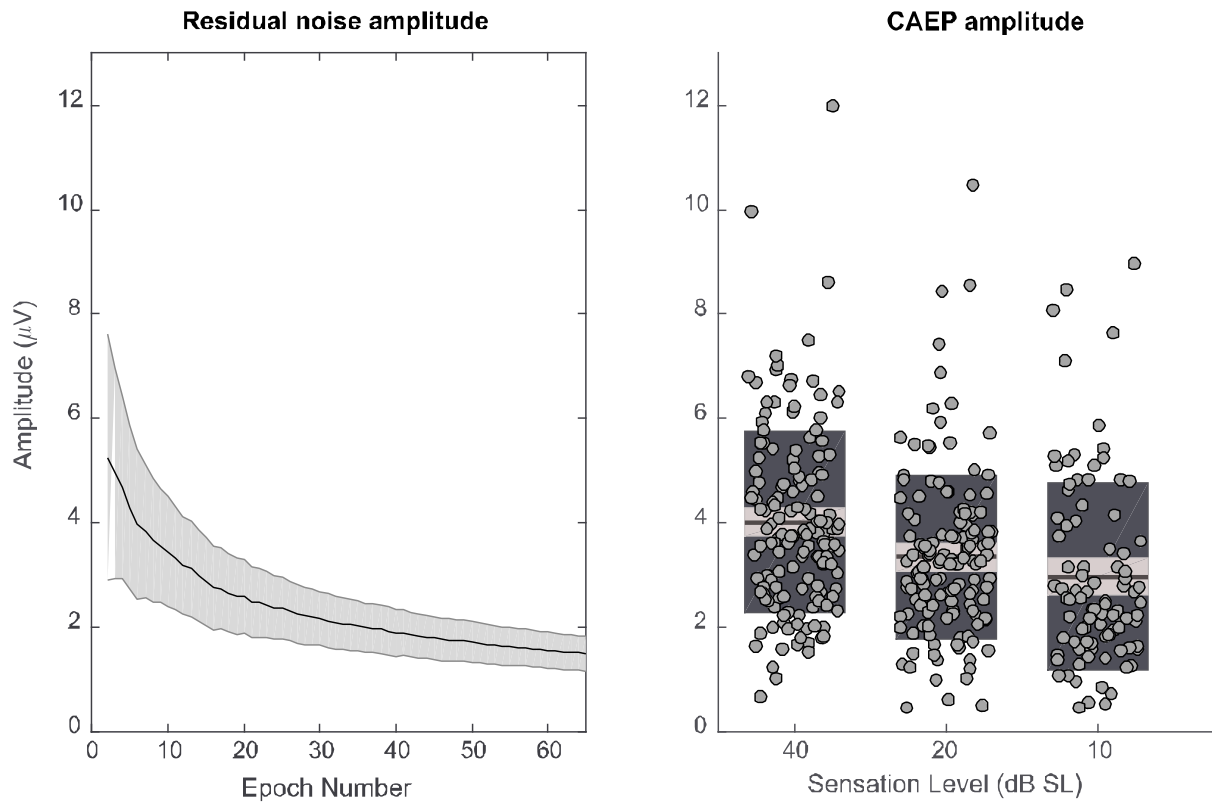
566

567

568

569 **Figure 1**

570

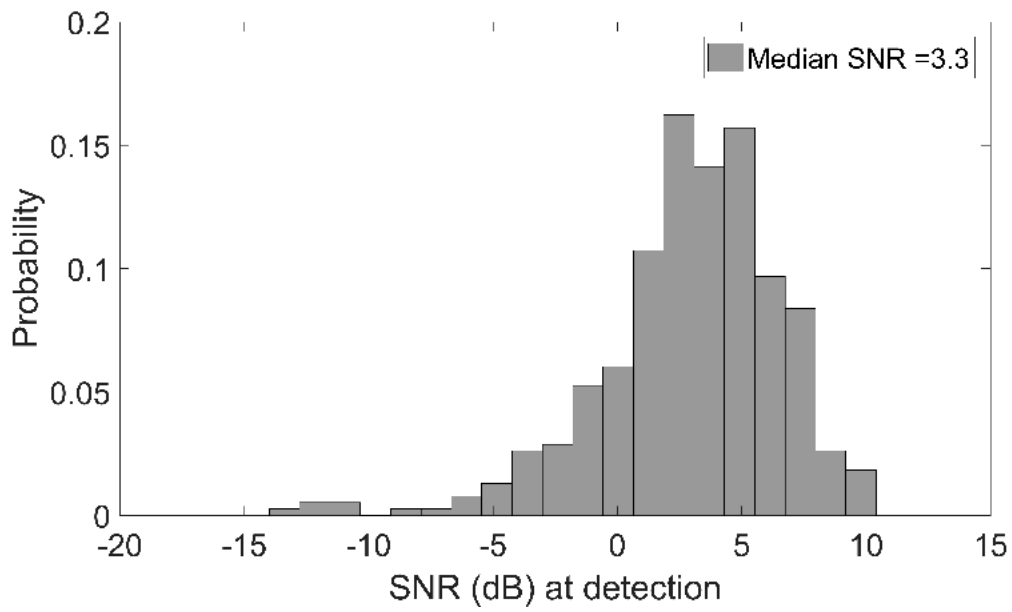


571

572 Figure 1. a) Represents the residual noise present in the EEG recording as a function of the number  
573 of epochs collected in a normal-hearing adult population. The grey shaded area represents the  
574 median and the epoch-to-epoch standard deviation of the RN rms amplitudes,. b) CAEP rms  
575 amplitude in normal-hearing adults plotted as a function of sensation levels (i.e. 40, 20 and 10 dB  
576 SL). For the signal and noise amplitudes to be comparable, residual noise was expressed as its rms  
577 amplitude, while the signal amplitude was expressed as as the square root of the difference  
578 between the average waveform power in the time window from 51 to 347 ms post-stimulus onset  
579 and the estimated residual noise power.

580

581 **Figure 2**

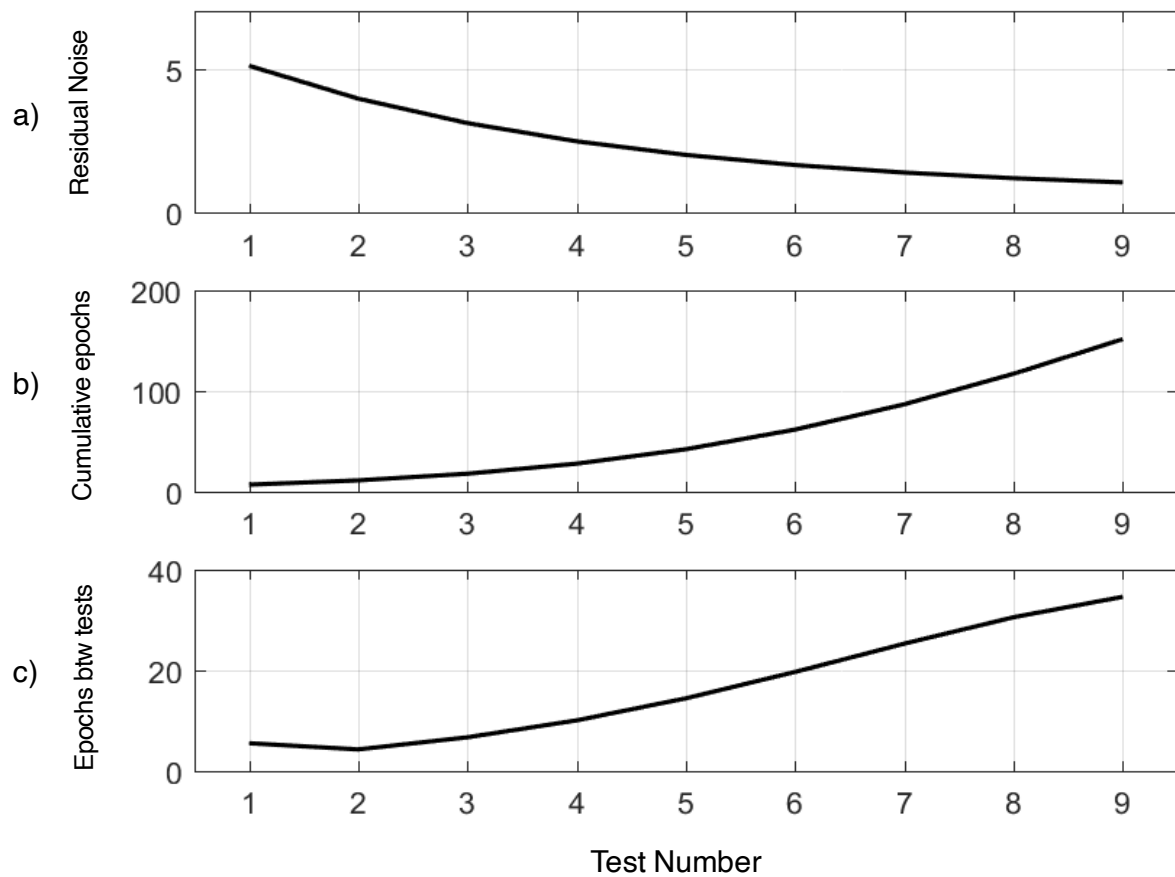


582

583 Figure 2: Distribution of SNRs when a CAEP is first detected using the Hotelling's T2 statistic  
584 constrained by a 5% FPR (normal-hearing adults).

585

586 **Figure 3**



587

588

589 Figure 3: a) Represents the RN amplitudes at which Hotelling's T2 statistical tests are conducted.

590 The equation of the RN criterion is  $6 \cdot \exp(-X/3.4) + 0.63$ . b) Number of epochs necessary to reach the

591 RN criteria displayed above for a subject having a residual noise level per epoch of  $12.5 \mu V$ . c)

592 Representation of the number of epochs between 2 tests before reaching the next noise criterion.

593

594 **Appendix**

595 We give an outline of the proof of the early stopping result. In the following, that result is called  
596 Corollary 2.

597 We use the same notation as in the body of the article:  $k$  is the number of bins,  $n$  is the "current"  
598 number of epochs,  $N$  is a number of epochs larger than  $n$ , and  $\Psi_{\nu_1, \nu_2}$  is the cumulative distribution  
599 function of an F random variable with degrees of freedom  $\nu_1$  and  $\nu_2$ .

600 The situation we consider is that we have observed  $n$  epochs  $x_1, \dots, x_n$ , with each  $x_i$  being a  $k$ -  
601 dimensional vector of numbers representing the  $i$ th epoch after binning. We then ask, if an additional  
602  $N - n$  epochs are added to the original sample to get an extended sample  $x_1, \dots, x_n, x_{n+1}, \dots, x_N$  of size  
603  $N$ , what values can the extended sample's Hotelling p-value take?

604 Let  $T_r^2$  and  $p_r$  be the Hotelling  $T^2$  statistic and the associated p-value, respectively, of the sample  
605  $x_1, \dots, x_r$ , for  $r = n$  and  $r = N$ .

606 We consider the original sample  $x_1, \dots, x_n$  to be fixed, so the values of  $T_n^2$  and  $p_n$  are fixed. We assume  
607  $T_n^2 > 0$ , which is equivalent to assuming  $p_n < 1$ , and we make the usual assumption that covariance  
608 matrices of the data are positive definite.

609 Under these assumptions, we have the following results.

610 Theorem 1. The maximum possible value of  $T_N^2$  is

611 
$$\left(\frac{N-1}{n}\right) \left( \left(\frac{N}{n-1}\right) T_n^2 + N - n \right).$$

612 Corollary 1. The minimum possible value of  $p_N$  is

613 
$$1 - \Psi_{k, N-k} \left[ \left(\frac{N-k}{n}\right) \left( \left(\frac{N}{n-k}\right) \Psi_{k, n-k}^{-1} (1 - p_n) + \frac{N-n}{k} \right) \right].$$

614 Corollary 2. If  $0 < q < 1$  and



615 
$$p_n > 1 - \Psi_{k,n-k} \left[ \left( \frac{n-k}{N} \right) \left( \left( \frac{n}{N-k} \right) \Psi_{k,N-k}^{-1}(1-q) - \frac{N-n}{k} \right) \right],$$

616 then  $p_N > q$ .

617 As an aside, it is easy to show that the minimum possible value of  $T_N^2$  is 0 and the maximum possible  
618 value of  $p_N$  is 1.

619 To prove Theorem 1, we start by considering the simplest case of  $N = n + 1$  and  $k = 1$ , that is, one  
620 additional epoch and one-dimensional data. Note that with one-dimensional data, the Hotelling  $T^2$   
621 statistic is the square of the one-sample  $t$ -test statistic.

622 For each  $r$ , let  $\bar{x}_r$ ,  $s_r^2$  and  $t_r^2$  be the sample mean, sample variance and squared  $t$  statistic, respectively,  
623 of the first  $r$  epochs. Note that the assumption that  $t_n^2$  is non-zero means that  $\bar{x}_n$  is non-zero.

624 We want an expression for  $t_{n+1}^2$  in which the only variable quantity is  $x_{n+1}$ , keeping in mind that the  
625 first  $n$  epochs are assumed to be fixed. By definition we have

626 
$$t_{n+1}^2 = \frac{(n+1)\bar{x}_{n+1}^2}{s_{n+1}^2}, \quad (A1)$$

627 so we want to express  $\bar{x}_{n+1}$  and  $s_{n+1}^2$  in terms of fixed quantities and  $x_{n+1}$ . The required expressions  
628 are

629 
$$\bar{x}_{n+1} = \frac{n\bar{x}_n + x_{n+1}}{n+1}$$

630 and

631 
$$s_{n+1}^2 = \frac{(n+1)(n-1)s_n^2 + n(\bar{x}_n - x_{n+1})^2}{n(n+1)}.$$

632 Substituting these into (A1) gives

633 
$$t_{n+1}^2 = \frac{n(n\bar{x}_n + x_{n+1})^2}{(n+1)(n-1)s_n^2 + n(\bar{x}_n - x_{n+1})^2}, \quad (A2)$$

634 and we define  $g(x_{n+1})$  to be the right-hand side of (A2), so  $t_{n+1}^2 = g(x_{n+1})$ .

635 We want to maximise the function  $g$ , so we look for points at which its derivative is zero. The derivative  
 636 of  $g$  can be written as

$$637 \quad g'(x_{n+1}) = \frac{2n(n+1)(n\bar{x}_n + x_{n+1})((n-1)s_n^2 + n\bar{x}_n^2 - n\bar{x}_n x_{n+1})}{[(n+1)(n-1)s_n^2 + n(\bar{x}_n - x_{n+1})^2]^2},$$

638 so the solutions of  $g'(x_{n+1}) = 0$  are  $x_{n+1} = -n\bar{x}_n$  and

$$639 \quad x_{n+1} = \bar{x}_n + \frac{(n-1)s_n^2}{n\bar{x}_n}. \quad (A3)$$

640 At the first of these solutions  $g$  is 0, and since  $g$  can't be negative, this must be a minimum. It can be  
 641 verified that at the second solution, the second derivative of  $g$  is negative, so that is where  $g$  is  
 642 maximum.

643 The maximum possible value of  $t_{n+1}^2$  therefore occurs when  $x_{n+1}$  has the value in (A3). Substituting  
 644 this value into (A2), and denoting the maximum possible value of  $t_{n+1}^2$  by  $\max t_{n+1}^2$ , gives

$$645 \quad \max t_{n+1}^2 = \left(\frac{n+1}{n-1}\right) t_n^2 + 1, \quad (A4)$$

646 which is Theorem 1 for  $N = n + 1$  and  $k = 1$ .

647 We now consider the case with  $N = n + 1$  and  $k > 1$ . Using the notation  $t^2(v_1, \dots, v_r)$  to mean the  
 648 squared one-sample  $t$  statistic of the univariate sample  $v_1, \dots, v_r$ , it is well-known (e.g., Johnson &  
 649 Wichern, 2007) that

$$650 \quad T_r^2 = \max_{a \neq 0} t^2(a'x_1, \dots, a'x_r),$$

651 with  $a$  and the  $x_i$  being viewed as  $k \times 1$  matrices and  $a'$  denoting the transpose of  $a$ .

652 Using this fact together with (A4) gives

$$653 \quad \max T_{n+1}^2 = \left(\frac{n+1}{n-1}\right) T_n^2 + 1,$$

654 which is Theorem 1 for  $N = n + 1$ . From this, the result for general  $N > n$  can be proved by induction.

655 Corollary 1 can be obtained by using the fundamental relation

656 
$$p_r = 1 - \Psi_{k,r-k} \left( \frac{r-k}{k(r-1)} T_r^2 \right)$$

657 (Johnson & Wichern, 2007) for  $r = n$  and  $r = N$  together with Theorem 1, and by noting that the p-  
658 value is minimised when the corresponding  $T^2$  is maximised.

659 Corollary 2 can be obtained by rearranging Corollary 1.

660

661 **Reference:**

662 Johnson, R. A., & Wichern, D. W. (2007). Applied multivariate statistical analysis (6th ed.). Upper

663 Saddle River, NJ: Pearson Prentice Hall.

664