



This is the author's version of a work that was accepted for publication in the following source:

Mao, D., H. Innes-Brown, M. A. Petoe, Y. T. Wong, and C. M. McKay. 2018. Cortical auditory evoked potential time-frequency growth functions for fully objective hearing threshold estimation. *Hearing Research*. **370**: 74-83.

doi: [10.1016/j.heares.2018.09.006](https://doi.org/10.1016/j.heares.2018.09.006)

**Notice:** Changes introduced as a result of publishing processes such as copy-editing and formatting may not be reflected in this document. For a definitive version of this work, please refer to the published source.

The final publication is available [here](#)

Copyright of this article belongs to: © 2018 Elsevier B.V.

# Accepted Manuscript

Cortical auditory evoked potential time-frequency growth functions for fully objective hearing threshold estimation

Darren Mao, Hamish Innes-Brown, Matthew A. Petoe, Yan T. Wong, Colette M. McKay



PII: S0378-5955(18)30252-1

DOI: [10.1016/j.heares.2018.09.006](https://doi.org/10.1016/j.heares.2018.09.006)

Reference: HEARES 7619

To appear in: *Hearing Research*

Received Date: 8 June 2018

Revised Date: 24 August 2018

Accepted Date: 26 September 2018

Please cite this article as: Mao, D., Innes-Brown, H., Petoe, M.A., Wong, Y.T., McKay, C.M., Cortical auditory evoked potential time-frequency growth functions for fully objective hearing threshold estimation, *Hearing Research* (2018), doi: <https://doi.org/10.1016/j.heares.2018.09.006>.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 Cortical auditory evoked potential time-frequency growth functions  
2 for fully objective hearing threshold estimation

3

4 Darren Mao<sup>a,b</sup>, Hamish Innes-Brown<sup>b,c</sup>, Matthew A. Petoe<sup>b,c</sup>, Yan T. Wong<sup>a,d</sup>, Colette M.  
5 McKay<sup>b,c</sup>

6 <sup>a</sup> Department of Biomedical Engineering, University of Melbourne, Parkville VIC 3010,  
7 Australia

8

9 <sup>b</sup> The Bionics Institute, 384-388 Albert St, East Melbourne VIC 3002, Australia

10

11 <sup>c</sup> Department of Medical Bionics, University of Melbourne, Parkville, VIC 3010, Australia

12

13 <sup>d</sup> Department of Physiology, Department of Electrical and Computer Systems Engineering,  
14 and the Biomedicine Discovery Institute, Monash University, Clayton VIC 3168, Australia

15

16 **Corresponding author:** Darren Mao, darren.mao92@gmail.com, Mailing address: The  
17 Bionics Institute, 384-488 Albert St, East Melbourne, VIC 3002, Australia

18 **Abstract**

19 Cortical auditory evoked potential (CAEPs) thresholds have been shown to correlate well  
20 with behaviourally determined hearing thresholds. Growth functions of CAEPs show promise  
21 as an alternative to single level detection for objective hearing threshold estimation; however,  
22 the accuracy and clinical relevance of this method is not well examined.

23 In this study, we used temporal and spectral CAEP features to generate feature growth  
24 functions. Spectral features may be more robust than traditional peak-picking methods where  
25 CAEP morphology is variable, such as in children or hearing device users. Behavioural  
26 hearing thresholds were obtained and CAEPs were recorded in response to a 1 kHz puretone  
27 from twenty adults with no hearing loss. Four features, peak-to-peak amplitude, root-mean-  
28 square, peak spectral power and peak phase-locking value (PLV) were extracted from the  
29 CAEPs. Functions relating each feature with stimulus level were used to calculate objective  
30 hearing threshold estimates. We assessed the performance of each feature by calculating the  
31 difference between the objective estimate and the behaviourally-determined threshold.

32 We compared the accuracy of the estimates using each feature and found that the peak PLV  
33 feature performed best, with a mean threshold error of 2.7 dB and standard deviation of 5.9  
34 dB across subjects from behavioural threshold. We also examined the relation between  
35 recording time, data quality and threshold estimate errors, and found that on average for a  
36 single threshold, 12.7 minutes of recording was needed for a 95% confidence that the  
37 threshold estimate was within 20 dB of the behavioural threshold, using the peak-to-peak  
38 amplitude feature, while 14 minutes is needed for the peak PLV feature. These results show  
39 that the PLV of CAEPs can be used to find a clinically relevant hearing threshold estimate.  
40 Its potential stability in differing morphology may be an advantage in testing infants or  
41 cochlear implant users.

42

43 **Keywords:** Hearing threshold; objective audiometry; electroencephalography; spectral  
44 analysis; phase-locking value; growth functions;

45

## 46 1 Introduction

47 Accurate estimations of hearing thresholds are important for the optimal adjustment of  
48 hearing aids and cochlear implants to suit each individual patient's hearing characteristics.  
49 Behavioural threshold estimates are the current gold standard and in most cases are easy to  
50 administer, but the accuracy of the result can be impaired by an inability to respond (as in the  
51 very young), lack of co-operation, or confusion on the part of the respondee. For these  
52 reasons, objective measures of hearing threshold are desirable. In this work, we use the  
53 cortical auditory evoked potential recorded from a single active channel as an objective  
54 measure used to estimate hearing thresholds, and compare the performance of different  
55 features of the cortical auditory evoked potential in threshold estimation using a growth  
56 function fitting method.

57 Objective hearing threshold estimates may be derived from auditory-evoked  
58 electroencephalography (EEG) measures such as the compound action potential (CAP)  
59 (Goldstein Jr et al., 1958), the auditory brainstem response (ABR) (Jewett et al., 1970), the  
60 auditory steady state response (ASSR) (Galambos et al., 1981), or the cortical auditory  
61 evoked potential (CAEP) (Davis, 1965). Electrophysiological recordings have the advantage  
62 that they do not require the participant to give feedback. When assessing the validity of these  
63 methods against behavioural thresholds, the CAEP, which is more likely than subcortical  
64 potentials to correspond to perception of the stimulus, appears to correlate with behavioral  
65 thresholds similarly or better than the ABR (Hoth, 1993; Pratt et al., 1978; Van Dun et al.,  
66 2015) and the ASSR (Tomlin et al., 2006; Yeung et al., 2007). The CAP (generated in the  
67 auditory nerve) is only utilised for cochlear implant users, where it is better known as the  
68 electrically-evoked CAP (ECAP), and is less correlated with behavioral hearing thresholds  
69 than either CAEP or ABR (Brown et al., 2000; McKay et al., 2017; McKay et al., 2013;  
70 Miller et al., 2008). In addition, the CAEP have shown promise in situations where ABR  
71 thresholds cannot be obtained (e.g. those with auditory neuropathy (Pearce et al., 2007; Starr  
72 et al., 1996)), or where ABR thresholds are not well correlated with behavioural measures  
73 (e.g. cochlear implant users, especially at higher stimulation rates (Abbas et al., 1991; Brown  
74 et al., 1999)).

75 Methods of hearing threshold estimation based on physiological responses use a common  
76 framework regardless of whether the estimation technique is by visual inspection or is  
77 automated by an algorithm. Estimation methods can generally be split into three steps: First,

78 multiple recordings are acquired in response to each of several sound intensity levels.  
79 Second, specific features are extracted from the EEG recording. Examples of these features  
80 include the frequently used N1-P2 potential difference (Picton et al., 1976), the phase  
81 coherence between individual response epochs of the same stimulus (Picton et al., 1987;  
82 Tallon-Baudry et al., 1996), the spectral power of the ASSR (Yeung et al., 2007), or the  
83 global field power across a full-scalp electrode montage (Visram et al., 2015). Third and  
84 finally, a set of decision rules are applied to find the threshold by either statistically testing  
85 for the presence of the feature at each single intensity level (Golding et al., 2009), or by  
86 fitting a growth function to extrapolate beyond the intensity levels that were tested to find  
87 threshold (Visram et al., 2015). Single intensity level testing requires using an adaptive  
88 procedure to find the lowest intensity level that elicits a response (Van Dun et al., 2015),  
89 however there is no literature on the optimal intensity levels to present for growth function  
90 fitting.

91 Previous studies have explored many combinations of the feature extraction and threshold  
92 estimation steps described above. By far the most commonly used feature in clinical use for  
93 threshold estimation is visual inspection of the N1-P2 morphology. These studies have  
94 consistently found a 10 dB standard deviation in difference between estimated threshold and  
95 behavioural threshold for normal hearing and hearing impaired adults using acoustic stimuli,  
96 although all have isolated exceptions where the estimate disagreed with behavioural estimates  
97 by over 15 dB or where no response could be recorded (Coles et al., 1984; Hoth, 1993;  
98 Rickards, 1996; Yeung et al., 2007). In objective detection, the amplitude of the time-domain  
99 response remains the most straightforward and most used feature (Hoth, 1993; Ross et al.,  
100 1999; Van Dun et al., 2015). However, other features such as single epoch phase (Ross et al.,  
101 1999), signal to noise ratio, reproducibility, and spectral power (Hoth, 1993) have been  
102 combined with amplitude measures, achieving results similar to those of visual detection by  
103 experienced examiners.

104 To further advance the robustness of objective threshold estimation, features based on the  
105 time-frequency content of the CAEP are of particular interest, as there has been no direct  
106 comparison of their performance against time-domain analyses. In addition, although power  
107 change from baseline (Hoth, 1993; Makeig, 1993) and phase measures (Ross et al., 1999)  
108 have been previously used in CAEP based analyses, the relationship of these features with  
109 stimulus intensity is largely unexplored. For example, Ross et al. (1999) reported an

110 increased statistical significance of phase locking compared to baseline phase locking when  
111 stimulus intensity was increased, but did not quantify the degree of phase locking.

112 Threshold estimation can benefit from utilising the relationship between features and  
113 stimulus intensity. This relationship is found by fitting growth functions to observed features  
114 of EEG recordings in response to sound stimuli of different intensities (see Figure 1). In  
115 comparison to methods that use single level statistical testing (e.g. Van Dun et al. (2015)) the  
116 growth function technique has the advantage that extrapolation of the growth function can  
117 give finer resolution of threshold estimates, and each tested level does not have to reach  
118 statistical significance. Visram et al. (2015) have used a growth function method successfully  
119 in cochlear implant users; however, 64 EEG channels were needed to estimate global field  
120 power. Use of fewer or single recording channel pairs may allow their direct incorporation  
121 onto hearing aids or cochlear implants, and hence audiometry in aided hearing without  
122 additional recording equipment. Growth functions of CAEPs are otherwise not widely  
123 reported due to the the time consuming nature of acquiring the necessary responses (Abbas et  
124 al., 2015).

125 In the present study, we explored time-frequency features of the EEG acquired from single  
126 EEG channels and their associated growth with stimulus intensity in order to estimate hearing  
127 thresholds in 20 normally hearing adults. We compared two time-domain and two time-  
128 frequency domain features: in the time domain, the peak-to-peak amplitude (N1-P2) and the  
129 total power in the response (root-mean-square value); and in the frequency domain, the post-  
130 stimulus peak spectral power and peak phase-locking value. We objectively estimated  
131 hearing thresholds with each of these features. We compared the performance of these four  
132 features in threshold estimate accuracy and established the relation between threshold  
133 estimate accuracy and recording time. We hypothesised that growth function fits based on  
134 time-frequency features would provide more accurate estimates of hearing threshold than  
135 growth function fits based on time domain features.

## 136 2 Methods

137 <<Insert Figure 1 about here, grayscale>>

## 138 2.1 *Subjects*

139 20 subjects, with ages ranging from 21 to 37 years (mean 25 years), participated in this study.  
140 This study was conducted under the approval of the Human Research Ethics Committee of  
141 the Royal Victorian Eye and Ear Hospital, and each participant gave written informed  
142 consent. An audiogram was performed for both ears at octave frequencies between 250 and  
143 8000 Hz with an audiometer and all subjects had thresholds less than 20 dB HL. The ear with  
144 the lower pure tone average threshold between 250 and 4000 Hz was selected for subsequent  
145 testing, resulting in 12 right ears and eight left ears.

## 146 2.2 *Sound stimuli*

147 The test stimulus was a 1000 Hz pure tone of 50ms duration and had 4-ms linear rising and  
148 falling edges, and phases were consistent for each stimulus. The stimuli were presented  
149 through Matlab (version 2015a, MathWorks, MA, USA), an RME Fireface 400 sound card  
150 (RME Audio, Germany) and Etymotic Research ER2 insert earphones (Etymotic Research,  
151 IL, USA). Stimulus calibration was performed with the GRAS 43AC Ear Simulator Kit and  
152 the Norsonic Nor140 Sound Analyser by taking a reading at a 1/3 octave band around a 1  
153 kHz centre frequency. The dB SPL of a full-scale pure tone was recorded and used to scale  
154 stimuli to the required dB SPL levels in Matlab.

## 155 2.3 *Behavioural threshold estimate*

156 Behavioural thresholds for the experimental stimulus were found with an adaptive three-  
157 interval forced-choice task. One of three 200-ms intervals, chosen randomly, contained the  
158 sound stimulus, where a 150ms silent period was appended to the end of the 50ms tone, while  
159 the other two contained silence. Subjects were asked to select the interval with sound or to  
160 guess if unsure. Stimuli began at 50 dB SPL, and the level was increased after one incorrect  
161 response or decreased after two consecutive correct responses. There were ten turning points,  
162 the first two being in steps of 10 dB and 5 dB, followed by eight 2 dB steps. The behavioural  
163 threshold was taken as the average of the last six turning points and is referred to as 0 dB  
164 sensation level (dB SL) in this paper.

## 165 2.4 *EEG data*

166 Nine stimulus intensities relative to 0 dB SL were presented: -20, -5, 0, 2, 5, 10, 20, 40 and  
167 60 dB SL. The -5, 0, and 2 dB SL levels were excluded from growth function analyses as  
168 they were deemed to not contribute to the growth function fits; instead they were used in

169 single level statistical testing. Additionally, the -20 dB SL level was only used to calculate  
170 baseline levels and was not used in the growth function fitting. To minimize possible effects  
171 of training or fatigue, stimuli were presented in a random order with an inter-stimulus interval  
172 (ISI) randomly jittered between 1.35 and 1.65 seconds. EEG data were acquired with a  
173 BioSemi ActiveTwo EEG system with data acquired from Cz (in the international 10-20  
174 configuration) and left and right mastoids, with the Common Mode Sense (CMS) and Driver  
175 Right Leg (DRL) placed on either side of the midpoint between Pz and POz. Data were  
176 sampled at 2048 Hz, with a fifth order sinc response anti-aliasing filter with -3dB cut-off at  
177 400 Hz. Data were acquired at each stimulus intensity until the residual noise level (Elberling  
178 et al., 1984) either reached  $0.4 \mu V^2$ , or  $0.6 \mu V^2$  if the Hotelling  $T^2$  statistic (Golding et al.,  
179 2009; Hotelling, 1992) concomitantly returned  $p < 10^{-6}$ . This ensured that the data quality  
180 across all subjects were at a consistent noise level. Recording was performed in an  
181 electrically isolated acoustic chamber. The BioSemi A/D box was connected to computer  
182 equipment outside the room via an optic cable. Subjects were given a silent, captioned film to  
183 watch during the experiment, and were instructed to stay awake and as still as possible. The  
184 noise from the stimulus generation system was measured to be 6.0 dB SPL.

### 185 2.5 *Data post-processing*

186 Data were re-referenced to the average of the two mastoid channels, and zero-phase filtered  
187 between 1 and 45 Hz with separate low and high-pass elliptical filters designed to have a  
188 passband ripple of  $<1\%$  and stopband attenuation of  $>40\text{dB}$ . Data were then down-sampled to  
189 256 Hz and epoched in a window -800 ms to +1200 ms relative to stimulus onset. Finally,  
190 any epochs with amplitudes that exceeded  $\pm 100 \mu V$  were rejected as artefact contaminated  
191 epochs. After these analyses, the dataset included data from 20 subjects and nine stimulus  
192 levels. We will henceforth refer to each set of epochs from one subject and one stimulation  
193 level as a “block” of epochs. The number of epochs ranged from 130 to 451 per block (mean:  
194 265, standard deviation: 86).

### 195 2.6 *Feature extraction*

196 Four response features of the EEG were considered in this study; two in the time domain, and  
197 two in the time-frequency domain. The two time domain features were both calculated on the  
198 epoch block mean (average of all epochs across each time point in a block): peak-to-peak  
199 amplitude values were calculated as the difference between the minimum and maximum  
200 voltage potential in a 50 to 500 ms post-stimulus window, and the root mean square (RMS)

201 value of the block mean was calculated on the same 50 to 500 ms post-stimulus window  
 202 (Figure 2A). The two time-frequency domain features, peak spectral power and peak phase  
 203 locking value (PLV), were calculated by first taking the short-time Fourier transform using  
 204 Matlab's spectrogram function on each epoch. A 400 ms Hamming window with a time-step  
 205 of 20 ms was used. Power spectrograms were calculated by taking ten times the log of each  
 206 epoch's time-frequency magnitude, averaging across all epochs, and then normalizing the  
 207 spectrogram by subtracting the pre-stimulus power calculated at 300 ms before stimulus  
 208 onset. The peak spectral power feature was defined as the peak value in the power  
 209 spectrogram between 1-20 Hz and between 50 and 500 ms, where a post-stimulus time is the  
 210 location of the centre of the Hamming window (Figure 2B). The phase-locking spectrograms  
 211 were calculated by taking each of the N epoch's time-frequency phase ( $\theta$ ) and applying the  
 212 following formula to calculate the phase-locking value (PLV) at each time ( $t$ ) and frequency  
 213 ( $f$ ) point to generate the phase-locking spectrogram (Mardia, 2014):

$$PLV(t,f) = \frac{1}{N} \sqrt{\left[ \sum_{i=1}^N \cos(\theta_i(t,f)) \right]^2 + \left[ \sum_{i=1}^N \sin(\theta_i(t,f)) \right]^2}$$

214 The peak PLV feature was the peak value in the PLV spectrogram between 1-20 Hz and  
 215 between 50 and 500 ms (Figure 2C).

## 216 2.7 Growth function fitting and threshold estimation

217 Each time a feature was extracted from a block of epochs in response to one stimulus  
 218 intensity, bootstrapping was performed to establish a confidence interval for each calculated  
 219 feature. The feature's distribution was estimated by randomly selecting a subset of epochs  
 220 from a block with replacement for 1000 iterations (Efron et al., 1986). The standard deviation  
 221 of the feature distribution was calculated (Efron et al., 1986), which gives a measure of the  
 222 size of the confidence interval. We will refer to this as the 'feature noise level'. Thus, a larger  
 223 feature noise level corresponded to a bigger uncertainty of the calculated feature. If a block  
 224 did not have enough epochs to reach a specified feature noise level, all epochs within the  
 225 block were used.

226 Single level testing was performed with the bootstrap feature distributions as above, and for  
 227 each subject. The bootstrap distributions for each level were compared to the baseline level.  
 228 The baseline level was the calculated feature using the -20 dB SL response block. Responses

229 at each level were deemed detected when the distributions of the tested level and baseline  
230 level overlapped for less than 50 bootstrap iterations (thus calculated  $p < 0.05$ ). The threshold  
231 estimate was defined as the lowest stimulus level at which two consecutive stimulus levels  
232 were detected as responses.

233 Growth function estimation was performed by finding the best linear fit that described the  
234 growth of the function with stimulus level, and finding the stimulus intensity level at which  
235 this function intersected with the feature noise floor. The median of the bootstrap feature  
236 distribution was used as the estimate of the feature at each stimulus intensity level. Growth  
237 functions were generated by performing a linear least-squares curve fit to each of the four  
238 response features using the 5, 10, 20, 40 and 60 dB SL intensity levels. Then, the feature  
239 value at -20 dB SL (an intensity level which would be very unlikely to be perceived) was used  
240 as the baseline value of this particular feature, and the dB SL level that corresponded to this  
241 value on the linear fit growth function was taken as the threshold estimate (see Figure 4A for  
242 diagrammatic representation). Thus, the threshold estimates described in this paper are the  
243 difference between the estimated threshold and the behavioural threshold in decibels.  
244 Threshold estimates were deemed invalid, and thus removed from further analyses, if one of  
245 three criteria were met: the estimated threshold was outside the range of  $\pm 100$  dB SL, the  
246 software was unable to find a valid fit, or the slope of the fit was less than 0 (negative  
247 amplitude growth). We also compared the linear curve fit to an exponential curve fit as  
248 described by Visram et al. (2015) and (Ross et al., 1999), fitting to the equation:

$$f(x) = a(1 - e^{-\frac{x-b}{c}})$$

249 where  $f$  represents the feature, and  $x$  represents the stimulus intensity in dB SL.  $a$  represents  
250 the asymptote and thus the upper limit for the features,  $b$  represents the x-axis shift and  $c$   
251 represents the speed of feature growth. When fitting a curve for the peak PLV feature,  $a$  is  
252 restricted to a value between 0 and 1.

253 Here, the threshold estimates were deemed invalid if the estimated threshold was outside the  
254 range of  $\pm 100$  dB SL, the adjusted  $r^2$  value of the curve fit was less than 0, or the term  $a$  or  $c$   
255 of the curve fit was less than 0. The arbitrary range of  $\pm 100$  dB SL was chosen to eliminate  
256 any outliers due to unusual fits (such as one with very shallow slope) affecting the overall  
257 analysis.

258 The performance of the threshold estimate method was quantified by the estimate's standard  
259 deviation across subjects. This method of defining performance was chosen, rather than the  
260 mean difference between behavioral and objective thresholds, because any mean difference  
261 can be subtracted from the objective threshold as a corrective factor (Hoth, 1993; Ross et al.,  
262 1999). For example, if a method estimated thresholds to be  $5 \pm 10$  dB SL across subjects,  
263 then we could subtract 5 dB from all estimates and the 10 dB standard deviation quantifies  
264 the error across subjects in the performance of this method.

### 265 3 Results

#### 266 3.1 Behavioural data

267 Hearing thresholds at 1 kHz determined from the audiogram across 20 subjects were  $13.0 \pm$   
268  $5.1$  dB SPL (mean  $\pm$  1 std). Behavioural thresholds obtained by the three-interval forced-  
269 choice task with 1 kHz tone pips were found to be  $14 \pm 4.4$  dB SPL (mean  $\pm$  1 std). A paired  
270 t-test revealed no significant difference between thresholds found with the two methods ( $t =$   
271  $1.57$ ,  $p = .133$ ), and there was a moderate and significant correlation between the methods ( $r$   
272  $= .76$ ,  $p < .001$ ). The three-interval forced-choice behavioural thresholds were taken as 0 dB  
273 SL and thus were used to compare with the objective estimates in further analyses.

#### 274 3.2 Significant power change and phase locking to sound stimuli were present in all 275 subjects

276 An example of epoch-averaged time-domain data from one subject at 60 dB SL is shown in  
277 Figure 2A, which illustrates how the two time-domain features, peak-to-peak amplitude and  
278 RMS were calculated. The same response represented in the time-frequency domain (spectral  
279 power and PLV) is shown in panels B and C, respectively. The peak values in each (indicated  
280 by a black dot) are the time-frequency domain features that were extracted for growth  
281 function fitting.

282 The regions enclosed by the black outline in Figures 2B and 2C are the regions within which  
283 the power and phase-locking features were statistically greater than baseline for that subject,  
284 with the baseline set at the -300 ms timepoint (permutation test, 1000 iterations, normalized  $z$   
285 score  $> 3$  for outline). The bottom panels (D-F) show the grand average time- and frequency-  
286 domain features, as well as the peak power and phase-locking locations for each of the twenty

287 subjects, denoted by black dots. It can be seen that the peak PLV time-frequency location was  
288 less variable across subjects than for peak power.

289 <<Insert figure 2 about here, in color>>

290 Each feature was extracted at every stimulus intensity for all subjects. The mean growth of  
291 the features with respect to stimulus intensity is shown in Figure 3. For all features, the values  
292 were observed to increase with increasing stimulus intensity. Furthermore, each feature in  
293 response to several near-threshold intensities were similar in value to the baseline feature  
294 level, where the baseline was indicated by the feature value at a stimulus intensity of -20 dB  
295 SL. Linear growth of each feature was observed at higher supra-threshold levels that were not  
296 obscured by the noise floor. This result prompted the selection of a linear curve fit for  
297 threshold estimation.

298 <<Insert figure 3 about here, in grayscale>>

### 299 3.3 *Linear curve fits estimate thresholds accurately*

300 Linear growth functions were fitted to each subject's data individually using the four features  
301 described in the Methods section. In this section, all epochs that were acquired were used in  
302 the analysis. In Figure 4, panel A shows an example of the growth function for one subject  
303 using the peak PLV feature. A linear relationship between the PLV and the stimulus intensity  
304 level can be observed. Panel B shows fitted linear growth functions for all subjects and  
305 features, and that peak PLV and peak-to-peak amplitude had the smallest variance across  
306 subjects in their threshold estimates, while peak spectral power was the worst performing  
307 measure with three invalid fits and one that was outside the displayed range.

308 The quality of the threshold estimation is shown in rows 1 and 2 of Table 1. Levene's test for  
309 unequal variance was performed in Minitab and showed a significant difference in the  
310 variance of the threshold estimates using the four features ( $W = 5.08, p = 0.003$ ). Post-hoc  
311 analysis by performing a two-sample  $F$ -test on each pairing found that the variance of the  
312 peak power threshold estimate was significantly higher than the variance of all the other  
313 features (Bonferroni corrected,  $\alpha = .05$ : versus peak-to-peak amplitude:  $F_{(16,19)} = 11.4, p <$   
314  $0.001$ ; versus RMS:  $F_{(16,19)} = 7.22, p < 0.001$  versus peak phase-locking value:  $F_{(16,19)} = 12.1,$   
315  $p < 0.001$ ). Thus, the peak power feature has poorer performance compared to the threshold  
316 estimates generated using the other three features. In addition, a one-way ANOVA was

317 performed on the mean threshold estimates and a significant difference was found (Welch's  
318 test,  $F_{(3,37)} = 3.43$ ,  $p = 0.027$ ). Post-hoc analysis revealed a significantly lower mean threshold  
319 estimate when using the peak PLV compared to the peak power feature (Games-Howell test,  
320 adjusted  $p = 0.029$ ). No other comparisons were significant.

321 <<Insert Figure 4 about here, in color>>

322 The mean goodness of fit across subjects for each feature, in adjusted  $r^2$  value, is indicated in  
323 row 2 of Table 1. The goodness of fit values were high (above 0.91 for all features except  
324 peak power), providing evidence that a linear function is suitable to model the relationship  
325 between the intensities tested and each of the features except for peak power. To further  
326 validate our method, we compared results in three other ways. Firstly, we compared the linear  
327 curve fit with an exponential curve fit (as performed by Visram et al. (2015)), and the  
328 goodness of fits are shown in Table 1. The threshold estimates and the goodness of fits are  
329 not significantly different for all four features (Ranksum test, Bonferroni corrected at  $\alpha =$   
330 0.05). However, by visual inspection, we observed that the exponential fits were mostly  
331 linear in the region near the hearing threshold and continued linearly for the higher stimulus  
332 intensities tested. The values for the upper asymptote of the fitted exponential for peak-to-  
333 peak amplitude, for example, were all above 90  $\mu\text{V}$  except for in one subject, an unlikely  
334 value for an auditory response. It was likely that the data did not sufficiently define the  
335 exponential fit, as the stimulus intensities did not reach a high enough level to for a plateau in  
336 the features to become apparent. In the special case of the peak PLV growth function where  
337 we limited the asymptote to a maximum of 1, the exponential growth functions could fit the  
338 sampled stimulus intensity points used in this study; and the growth rate was also slow  
339 enough to fit the mostly linear growth as seen in Figure 3 panel D. The implications of using  
340 different growth function shapes are included in the Discussion.

341 Secondly, instead of using the noise floor to estimate the threshold from the linear growth  
342 function, we used a value of zero. The resultant threshold estimates indicate that finding the  
343 intersection of the noise floor with the linear fit gave threshold estimates with lower standard  
344 deviations across subjects for peak to peak amplitude, RMS and peak PLV ( $F$ -test,  
345 Bonferroni corrected,  $\alpha = 0.05$ ;  $F_{19,19} = 3.83, 3.85, 6.87$  and  $p = 0.005, 0.005, 0.0001$   
346 respectively), while peak power did not ( $F_{16,16} = 2.71$ ,  $p = 0.054$ ).

347 Lastly, we compared linear curve fits to single level statistical testing. In Table 2, we  
 348 compared the Hotelling test, as used by (Van Dun et al., 2015) to the four features in this  
 349 paper using single level detection as described in the section 2.7 of Methods. A direct  
 350 comparison between thresholds estimated using a growth function and using single level  
 351 detection is difficult because different stimulus levels and data qualities at each level is  
 352 required for either method. This difficulty is further explained in the Discussion.  
 353 Nevertheless, we found no significant difference in the variance of the threshold estimates  
 354 using the five different forms of single level testing as described in the Methods section  
 355 (Levene's test,  $W = 0.809$ ,  $p = 0.523$ ).

356 **Table 1: Threshold estimates using all data for each of the features.** Mean  $\pm$  1 standard deviation and  
 357 goodness of fit (adjusted  $r^2$  value) shown. The first two rows show the results generated from the procedure as  
 358 described in the Methods section: linear curve fits, and threshold calculated using the feature value of baseline  
 359 response. The next two rows show results using an exponential curve fit, and the last row show results using the  
 360 linear curve fit but the baseline response is assumed to have a feature value of zero.

		Time-domain features		Time-frequency domain features	
		Amplitude	RMS	Power	PLV
Proposed Method	Threshold estimate (dB SL +- SD)	$5.7 \pm 5.9$	$5.4 \pm 7.6$	$18.5 \pm 20.2^*$	$2.7 \pm 5.9$
	Goodness of fit (adjusted $r^2$ )	$0.93 \pm 0.07$	$0.94 \pm 0.06$	$0.64 \pm 0.29^*$	$0.91 \pm 0.10$
Exponential fit	Threshold estimate (dB SL)	$5.4 \pm 5.8$	$6.5 \pm 5.4$	$10.4 \pm 10.0^\dagger$	$5.18 \pm 3.89$
	Goodness of fit (adjusted $r^2$ )	$0.89 \pm 0.17$	$0.91 \pm 0.09$	$0.61 \pm 0.34^\dagger$	$0.88 \pm 0.18$
Zero for noise floor	Threshold estimate (dB SL)	$-17.4 \pm 11.7$	$-21.4 \pm 14.8$	$-38.1 \pm 33.3^*$	$-21.1 \pm 15.3$

361 \*Three subjects excluded due to invalid fit (see text)

362 †Eight subjects excluded due to invalid fit

363 **Table 2:** Threshold estimates generated from single level statistical tests. The Hotelling test is as described in  
 364 (Van Dun et al., 2015), while the features explored in this study are tested for significance by comparing against  
 365 their bootstrapped baseline (-20 dB SL epochs). The lowest level at which a cortical response is significantly  
 366 different from baseline with the level immediately above it also detecting a response is deemed the threshold for  
 367 each of the five methods.

	<b>Hotelling</b>	<b>Amplitude</b>	<b>RMS</b>	<b>Power</b>	<b>PLV</b>
Threshold estimate (dB SL)	$17.8 \pm 10.8$	$26.7 \pm 15.9$	$29.5 \pm 16.7$	$48.3 \pm 15.9^\dagger$	$26.0 \pm 14.3$

368 <sup>†</sup>Eight subjects had responses that were not detected at any level

### 369 3.4 *Peak-to-peak amplitude measures perform better than other features when a limited* 370 *number of epochs are used*

371 The results in Figure 4 show excellent threshold estimates, particularly for the peak-to-peak  
 372 amplitude and peak PLV feature; however, up to 50 minutes' worth of data (mean = 38  
 373 minutes, std = 9.1 minutes) contributed to each threshold estimate. Therefore, we explored  
 374 the effect of reducing the number of epochs acquired (and hence the testing time needed) on  
 375 the threshold estimate accuracy. We expected that the feature noise level determines the  
 376 accuracy of the threshold estimate, where increased feature noise level corresponds to  
 377 decreased threshold estimate accuracy, since the feature noise level relates to the uncertainty  
 378 of the feature points used in the growth function fit. We also expected that the number of  
 379 epochs determines the feature noise level, where fewer epochs corresponds to increased  
 380 feature noise level, since we expect the uncertainty of the feature points to decrease as more  
 381 data is acquired. Thus, the expectation is that the threshold estimate accuracy decreases when  
 382 using fewer epochs. We quantify the relation between the number of epochs used for feature  
 383 calculations, the feature noise level and the threshold estimate accuracy in Figure 5.

384 <<Insert Figure 5 about here, in color>>

385 Figure 5A shows the relation between threshold estimate accuracy and feature noise level,  
 386 using the peak PLV feature as an example. 300 feature noise levels, evenly spaced between  
 387 0.03 and 0.1 (arbitrary units), were used as feature noise thresholds. For each of the stimulus  
 388 intensities, epochs were randomly sampled and added to the peak PLV calculation until the  
 389 feature noise level was below the feature noise threshold. If the feature noise threshold could  
 390 not be reached, all epochs in the block were used, and any requested feature noise levels

391 lower than the aforementioned were excluded from further analysis as they were duplicates  
392 using the same data. Following this, the growth function fitting and threshold estimation was  
393 performed. Thus, 6000 threshold estimates were calculated (300 feature noise levels for each  
394 of the 20 subjects) and each threshold estimate is represented by an individual point in Figure  
395 5A. We observed that the standard deviation across subjects of the threshold estimates,  
396 shown by the outer black lines, increases with higher feature noise level, indicating a  
397 decrease in threshold estimate accuracy.

398 Figure 5B shows the relation between feature noise level and the number of epochs used,  
399 using one subject's data as an example. We observed a clear, monotonic decrease in feature  
400 noise level as more epochs were added, regardless of stimulus intensity. In conjunction with  
401 Figure 5A, it can be inferred that threshold estimate accuracy decreases with fewer epochs,  
402 and accordingly, less recording time.

403 We quantified the relation between recording time and threshold estimate accuracy and the  
404 results are shown in panel C. A moving window of 30 feature noise levels was used to  
405 aggregate the threshold points for each feature (an example of these points estimated with  
406 peak PLV is shown in panel A), and the spread of these threshold estimates gave a  
407 quantifiable performance measure. For example, using the PLV feature guaranteed with 95%  
408 confidence that, at a feature noise threshold of between 0.063 and 0.07, the threshold estimate  
409 was within 20 dB of the actual threshold. The number of epochs that were used in generating  
410 each threshold estimate was then converted to a recording time by assuming a constant 1.5  
411 second ISI, resulting in  $14.0 \pm 1.7$  minutes of recording time to obtain the threshold estimate.  
412 Recording time estimates are summarised in Table 1. These results indicate that when using  
413 the peak-to-peak amplitude feature for example, we can be 95% confident that a threshold  
414 estimate is within 20 dB of the behavioural threshold after recording for an average of 12.7  
415 minutes.

416 These results are summarised in Table 3 by converting the feature noise level to recording  
417 time across all subjects. A two-way repeated measures ANOVA was performed in Minitab  
418 with the General Linear Model function with feature and threshold estimate accuracy  
419 restriction both as fixed factors. The model found significant differences in feature ( $F_{2,3594} =$   
420 1343), threshold estimate accuracy restriction ( $F_{1,3594} = 2819$ ), and the interaction between  
421 feature and threshold estimate accuracy restriction ( $F_{2,3594} = 175.2$ ), with  $p < 0.001$  for all three  
422 cases. Post-hoc analysis using the Tukey pairwise comparisons revealed that each of the peak

423 amplitude, RMS and peak PLV features took significantly different recording times to reach  
 424 the threshold estimate accuracy restriction ( $p < 0.001$  for all two feature comparisons across  
 425 two threshold estimate accuracy restrictions).

426 **Table 3.** Recording time needed to reach 95% confidence that one threshold estimate is within 15 dB and 20 dB  
 427 of the behavioural threshold. The results are expressed as the mean recording time in minutes with one standard  
 428 deviation across all simulated threshold estimates that was at a 15 dB SL and 20 dB SL 95% confidence interval  
 429 respectively.

Recording time needed for one threshold estimate to be:	Time domain features		Time-frequency domain features	
	Amplitude	RMS	Power	PLV
95% CI within $\pm 15$ dB (mean $\pm$ std. minutes)	18.0 $\pm$ 4.1	28.6 $\pm$ 6.0	Never reached	19.4 $\pm$ 2.3
95% CI within $\pm 20$ dB (mean $\pm$ std. minutes)	12.7 $\pm$ 3.1	18.0 $\pm$ 5.2		14.0 $\pm$ 1.7

430

#### 431 4 Discussion

432 The main aim of this work was to compare features extracted from CAEPs to acoustic stimuli  
 433 for use in hearing threshold estimation using growth function fitting. The significance of this  
 434 work lies in its clinical value, where objective and accurate hearing threshold estimates are  
 435 desired. Although we did not find a significant improvement by using time-frequency domain  
 436 features, we have shown that growth function fitting of root-mean-square value and peak  
 437 phase-locking value in addition to traditional amplitude growth functions are a viable method  
 438 for hearing threshold estimates. In addition, we have validated the linear relationship between  
 439 CAEP features and stimulus intensity in dB SL. The analyses defining minimum required  
 440 recording time is the first step towards a fully automated threshold estimation technique. We  
 441 have also shown that the accuracy of hearing threshold estimates is related to the recording  
 442 time, allowing comparison of results for hearing threshold estimation using novel growth  
 443 function algorithms.

#### 444 4.1 *Performance of features for threshold estimation*

445 We observed similar performance between peak-to-peak amplitude, RMS and peak PLV for  
446 threshold estimation, while the performance using spectral power is far worse than the other  
447 features. For the growth function fit to accurately estimate threshold, the features must be  
448 present with lower but still supra-threshold stimuli. However, it has been shown that there is  
449 no significant change in spectral power of the CAEP when the stimulus is at a low intensity  
450 (around 30 dB SL), whereas significant phase changes are present (Savers et al., 1974).

451 Furthermore, the similarity in performance between the other three features may be due to  
452 event-related dynamics represented by the features. As peak-to-peak amplitude and RMS  
453 operate on the averaged waveform, and peak PLV extracts phase locking of responses to  
454 individual trials, these methods extract responses conforming to the partial phase resetting  
455 (PPR) and the event-related potential (ERP) paradigms. As a change in spectral power from  
456 baseline would also be present in the ERP paradigm yet has poor performance in this study, it  
457 is possible that PPR is the underlying event-related dynamic being described by amplitude,  
458 RMS and PLV features. Harris et al. (2014) showed a correlation between P2 amplitude and  
459 PLV, but not spectral power, thus also suggesting the PPR paradigm, in their case for cortical  
460 responses to gaps in noise. The event-related dynamics represented by the different features  
461 could be further tested by classification using features. By classifying CAEPs from  
462 background brain activity using different combinations of features, a mutual information  
463 metric between actual response and predicted response could provide further insights into  
464 features extracted from EEG and their underlying event-related dynamics.

#### 465 4.2 *Growth of CAEP features with stimulus intensity*

466 We made two assumptions in our method for estimating hearing thresholds from growth  
467 functions. First, we assumed that the growth of these features with stimulus intensity is linear,  
468 at least in the intensity range of interest (threshold up to 60 dB SL). Although exponential fits  
469 have been used previously (Abbas et al., 2015; Ross et al., 1999; Visram et al., 2015), we  
470 have shown that the linear curves fit equally well. The differences in the types of functions  
471 used in these studies may be due to several factors. Firstly, there is a difference between  
472 acoustic and electric hearing, although these differences have not been quantified. One  
473 possibility is the much larger dynamic range in decibels of acoustic listeners than electric  
474 listeners; this may cause the growth functions in acoustic listeners to have a steeper slope or  
475 to change its shape at higher levels. Secondly, different EEG response features were used in

476 each: global field power gave a mix of linear and exponential fits across individual subjects  
477 (Visram et al., 2015) and N1 amplitude gave an exponential fit (Ross et al., 1999). Finally,  
478 Ross et al. (1999) used a grand average growth function to examine the characteristics of  
479 amplitude feature growth with stimulus intensity, but in our results (Figure 4B) and in those  
480 of Visram et al. (2015) and Abbas et al. (2015), individual growth functions differ. Using a  
481 grand average growth function may instigate the negative effects of inter-subject variability  
482 and give inaccurate growth functions and thus poorer threshold estimates.

483 It must also be noted that even though the growth function fit gives accurate hearing  
484 threshold estimates, the shape of the true underlying growth function is difficult to determine.  
485 A linear function would not be realistic at higher intensity levels because the features cannot  
486 continue growing; likewise, an exponential function would not be realistic at near- or below-  
487 threshold levels as the features are expected to be near the baseline value and cannot be  
488 negative. The accurate threshold estimates, however, suggest both functions are a good  
489 approximation of the shape of the underlying growth function, at least for stimulus levels  
490 within the person's dynamic range of hearing.

491 Second, we found that inferring behavioural threshold from the noise floor and fitted growth  
492 function gave better hearing threshold estimates than using a noise floor of zero (Table 1).  
493 This contrasts with both Visram et al. (2015) and Ross et al. (1999), who assumed that the  
494 threshold is the stimulus intensity at which the fitted function value is zero. Our results  
495 provide evidence that the CAEP features for the supra-threshold levels contain an amount of  
496 noise in addition to that of the underlying response. For the peak-to-peak amplitude, peak  
497 PLV and peak spectral power features, choosing the peak value introduces a positive bias to  
498 the feature; moreover, for the RMS and PLV, there is an inherent positive bias due to the way  
499 they are calculated. RMS includes the spontaneous EEG activity, while the baseline PLV is  
500 distributed based upon the number of epochs in its calculation. Performing the same  
501 calculations on the baseline level improves threshold estimation performance and provides  
502 evidence that it at least partially accounts for the aforementioned biases.

#### 503 4.3 *The effect of recording time on threshold estimation accuracy*

504 We have defined the notion of feature noise level, which describes the uncertainty in the  
505 features, becomes smaller as more data is acquired. We have subsequently shown that  
506 objective hearing threshold estimation accuracy is related to the feature noise level. This is  
507 useful knowledge for clinical applications as the trade-off between recording time and

508 accuracy can now be quantified, and a more informed choice can be made when designing a  
509 data acquisition paradigm.

510 Hoth (1993) and Van Dun et al. (2015) found a standard deviation of 11.4 and 10.2 dB  
511 difference from the behavioural threshold respectively, with a small but considerable number  
512 of outliers that had a difference of greater than 30 dB. In this study, we found a standard  
513 deviation of 5.9 dB using both the amplitude and the PLV feature with no differences greater  
514 than 20 dB. This came at the expense of recording for much longer than the previous studies.  
515 While Hoth (1993) did not give an indication of test time as more stimulus intensities than  
516 needed were presented, similar recording time per level was used compared to Van Dun et al.  
517 (2015) (50 sweeps with 2.5s ISI, 120 sweeps with 1.175s ISI respectively). Van Dun et al.  
518 used  $10.6 \pm 1.2$  minutes to reach standard deviation of the estimate of 10 dB; in this study, we  
519 achieved a similar result (95% within 20 dB, equal to a standard deviation of 10.2 dB) this in  
520  $12.7 \pm 3.1$  minutes using the peak-to-peak amplitude feature. Our study shows that the test  
521 time is slightly longer, meaning further improvements can consolidate growth function fitting  
522 as the preferred method for objective threshold estimation.

#### 523 4.4 *Comparison between growth functions and statistical testing for threshold estimates*

524 Growth function fitting is advantageous over statistical testing in two ways. Firstly, a set of  
525 measured responses which do not have significant amplitude change or phase locking does  
526 not imply that those epochs do not contain a response to the sound; rather, the model that has  
527 been built may not have sufficient epochs or statistical power to detect a significant response.  
528 The growth function fitting method can still make use of the information in epoch blocks in  
529 which the chosen feature is not statistically significant. Secondly, in single-level statistical  
530 testing, the accuracy of the threshold estimate is limited by the spacing of the stimulus levels  
531 tested. We have shown this in Table 2 where the hearing threshold estimates using single  
532 level testing are shown to have far greater variance across subjects than when using the  
533 growth function method. As previously mentioned, direct comparison between these methods  
534 is biased, as a true single-level detection algorithm would adaptively test levels (Van Dun et  
535 al. 2015) instead of larger, 20 dB steps.

536 The stopping criterion for acquiring epochs is a possible disadvantage of the growth function  
537 method compared to statistical testing. This is important as recording time is a key factor for  
538 whether a method is clinically viable. For statistical testing, false detection rates can be  
539 controlled, and acquisition of an epoch block can be terminated as soon as significance is

540 detected. In this work, we remedied this by showing that the feature noise level, and the  
541 number of epochs used, is a useful measure for the accuracy of the threshold estimate. The  
542 feature noise level can be used as a stopping criterion, whereas the number of epochs used  
543 gives an indication of how quickly the threshold estimate can be obtained.

#### 544 4.5 *Unexplored parameter space for growth function fitting*

545 There is a large, unexplored parameter space with the growth function method which may  
546 impact the recording time and accuracy of the threshold estimate. We will consider clinical  
547 relevance when discussing the following points. Firstly, the number of stimulus levels cannot  
548 be pre-determined in a clinical setting. For this method to be clinically applicable, an  
549 algorithm that will determine the best stimulus levels and dynamically test a range of levels is  
550 needed. Further work is needed on exploring the optimal number of points to fit the growth  
551 function: to optimise recording time, the relation between threshold estimate accuracy,  
552 number of stimulus levels and feature noise level needs to be identified. Secondly, the feature  
553 baselines were calculated from a well-below-threshold level in this study. It is possible to  
554 explore the use of the EEG response in the silent region between stimuli to calculate the  
555 baseline instead, saving further time. Lastly, audiograms or cochlear implant fitting generally  
556 require finding more than one hearing threshold at a time; testing multiple levels where  
557 different tones are mixed can shorten recording time further (Lightfoot et al., 2006).

#### 558 5 Conclusion

559 In this study, we have shown for the first time the performance of growth function fitting to  
560 objectively estimate hearing thresholds in normal hearing adults, using the CAEP recorded  
561 from one active electrode. The results show that linear growth functions produce accurate  
562 threshold estimates for all subjects with any of peak-to-peak amplitude, RMS or peak PLV  
563 features. In addition, growth function hearing threshold estimates are comparable to existing  
564 methods for fully objective hearing threshold estimate using the cortical auditory evoked  
565 response. Future work to explore other parameters of growth function fitting is necessary to  
566 further improve hearing threshold estimation.

## 567 6 Acknowledgements

568 Darren Mao was supported by an Australian Government Research Training Program Scholarship.

569 Colette McKay was supported by a Veski innovation fellowship. Hamish Innes-Brown was supported

570 by a National Health and Medical Research Council (Australia) Early-Career Research Fellowship.

571 The Bionics Institute acknowledges the support it receives from the Victorian Government through its

572 Operational Infrastructure Support Program.

573

## 574 7 References

- 575 Abbas, P.J., Brown, C.J. 1991. Electrically evoked auditory brainstem response: growth of  
576 response with current level. *Hear Res* 51, 123-37.
- 577 Abbas, P.J., Brown, C.J. 2015. Assessment of responses to cochlear implant stimulation at  
578 different levels of the auditory pathway. *Hear Res* 322, 67-76.
- 579 Brown, C.J., Lopez, S.M., Hughes, M.L., Abbas, P.J. 1999. Relationship between EABR  
580 thresholds and levels used to program the CLARION® speech processor. *Ann Otol*  
581 *Rhinol Laryngol* 108, 50-57.
- 582 Brown, C.J., Hughes, M.L., Luk, B., Abbas, P.J., Wolaver, A., Gervais, J. 2000. The  
583 relationship between EAP and EABR thresholds and levels used to program the  
584 nucleus 24 speech processor: Data from adults. *Ear Hear* 21, 151-163.
- 585 Coles, R.R., Mason, S.M. 1984. The results of cortical electric response audiometry in  
586 medico-legal investigations. *Br J Audiol* 18, 71-8.
- 587 Davis, H. 1965. Slow cortical responses evoked by acoustic stimuli. *Acta Oto-laryngologica*  
588 59, 179-185.
- 589 Efron, B., Tibshirani, R. 1986. Bootstrap methods for standard errors, confidence intervals,  
590 and other measures of statistical accuracy. *Stat Sci*, 54-75.
- 591 Elberling, C., Don, M. 1984. Quality estimation of averaged auditory brainstem responses.  
592 *Scand Audiol* 13, 187-97.
- 593 Galambos, R., Makeig, S., Talmachoff, P.J. 1981. A 40-Hz auditory potential recorded from  
594 the human scalp. *Proc Natl Acad Sci U S A* 78, 2643-7.
- 595 Golding, M., Dillon, H., Seymour, J., Carter, L. 2009. The detection of adult cortical auditory  
596 evoked potentials (CAEPs) using an automated statistic and visual detection. *Int J*  
597 *Audiol* 48, 833-42.
- 598 Goldstein Jr, M.H., Kiang, N.Y.S. 1958. Synchrony of neural activity in electric responses  
599 evoked by transient acoustic stimuli. *J Acoust Soc Am* 30, 107-114.
- 600 Harris, K.C., Vaden, K.I., Jr., Dubno, J.R. 2014. Auditory-evoked cortical activity:  
601 contribution of brain noise, phase locking, and spectral power. *J Basic Clin Physiol*  
602 *Pharmacol* 25, 277-84.
- 603 Hotelling, H. 1992. The generalization of Student's ratio, *Breakthroughs in Statistics*.  
604 Springer. pp. 54-65.
- 605 Hoth, S. 1993. Computer-aided hearing threshold determination from cortical auditory  
606 evoked potentials. *Scand Audiol* 22, 165-77.
- 607 Jewett, D.L., Romano, M.N., Williston, J.S. 1970. Human auditory evoked potentials:  
608 possible brain stem components detected on the scalp. *Science* 167, 1517-8.
- 609 Lightfoot, G., Kennedy, V. 2006. Cortical electric response audiometry hearing threshold  
610 estimation: accuracy, speed, and the effects of stimulus presentation features. *Ear*  
611 *Hear* 27, 443-56.
- 612 Makeig, S. 1993. Auditory event-related dynamics of the EEG spectrum and effects of  
613 exposure to tones. *Electroen Clin Neuro* 86, 283-93.
- 614 Mardia, K.V. 2014. *Statistics of directional data* Academic press.
- 615 McKay, C.M., Smale, N. 2017. The relation between ECAP measurements and the effect of  
616 rate on behavioral thresholds in cochlear implant users. *Hear Res* 346, 62-70.
- 617 McKay, C.M., Chandan, K., Akhoun, I., Siciliano, C., Kluk, K. 2013. Can ECAP Measures  
618 Be Used for Totally Objective Programming of Cochlear Implants? *J Assoc Res*  
619 *Otolaryngol* 14, 879-890.
- 620 Miller, C.A., Brown, C.J., Abbas, P.J., Chi, S.L. 2008. The clinical application of potentials  
621 evoked from the peripheral auditory system. *Hear Res* 242, 184-97.

- 622 Pearce, W., Golding, M., Dillon, H. 2007. Cortical auditory evoked potentials in the  
623 assessment of auditory neuropathy: two case studies. *J Am Acad Audiol* 18, 380-90.
- 624 Picton, T.W., Woods, D.L., Baribeau-Braun, J., Healey, T.M. 1976. Evoked potential  
625 audiometry. *J Otolaryngol* 6, 90-119.
- 626 Picton, T.W., Vajsar, J., Rodriguez, R., Campbell, K.B. 1987. Reliability estimates for  
627 steady-state evoked potentials. *Electroencephalogr Clin Neurophysiol* 68, 119-31.
- 628 Pratt, H., Sohmer, H. 1978. Comparison of hearing threshold determined by auditory pathway  
629 electric responses and by behavioural responses. *Audiology* 17, 285-92.
- 630 Rickards, F.W. 1996. Cortical evoked response audiometry in noise induced hearing loss  
631 claims. *Aust J Otolaryng* 2, 237.
- 632 Ross, B., Lutkenhoner, B., Pantev, C., Hoke, M. 1999. Frequency-specific threshold  
633 determination with the CERAGram method: basic principle and retrospective  
634 evaluation of data. *Audiol Neurootol* 4, 12-27.
- 635 Savers, B.M., Beagley, H., Henshall, W. 1974. The mechanism of auditory evoked EEG  
636 responses. *Nature* 247, 481-483.
- 637 Starr, A., Picton, T.W., Sininger, Y., Hood, L.J., Berlin, C.I. 1996. Auditory neuropathy.  
638 *Brain* 119 ( Pt 3), 741-53.
- 639 Tallon-Baudry, C., Bertrand, O., Delpuech, C., Pernier, J. 1996. Stimulus specificity of  
640 phase-locked and non-phase-locked 40 Hz visual responses in human. *J Neurosci* 16,  
641 4240-9.
- 642 Tomlin, D., Rance, G., Graydon, K., Tsialios, I. 2006. A comparison of 40 Hz auditory  
643 steady-state response (ASSR) and cortical auditory evoked potential (CAEP)  
644 thresholds in awake adult subjects. *Int J Audiol* 45, 580-588.
- 645 Van Dun, B., Dillon, H., Seeto, M. 2015. Estimating hearing thresholds in hearing-impaired  
646 adults through objective detection of cortical auditory evoked potentials. *J Am Acad*  
647 *Audiol* 26, 370-83.
- 648 Visram, A.S., Innes-Brown, H., El-Dereby, W., McKay, C.M. 2015. Cortical auditory evoked  
649 potentials as an objective measure of behavioral thresholds in cochlear implant users.  
650 *Hear Res* 327, 35-42.
- 651 Yeung, K.N., Wong, L.L. 2007. Prediction of hearing thresholds: Comparison of cortical  
652 evoked response audiometry and auditory steady state response audiometry  
653 techniques. *Int J Audiol* 46, 17-25.

**Figure 1: Representation of the threshold estimation method used in this study.** Multiple sound intensity levels are presented, labeled as  $L_1$ ,  $L_2$ ,  $L_3$  as well as a sub-threshold intensity level,  $L_{\text{sub}}$ ; multiple epochs of responses to each are recorded. Features are then extracted from each block of recorded EEG activity, labeled as  $f_{L_1}$ , etc. The baseline feature estimate is the feature extracted from baseline EEG activity, as there should be no response to a sub-threshold stimulus. For growth function threshold estimation, a linear curve is fit to all EEG features except the baseline feature estimate, and the sound stimulus intensity that corresponds to the sub-threshold response feature is the hearing threshold estimate.

**Figure 2: Top row: Responses to 60 dB SL stimuli for subject 14, represented in time and time-frequency domain, with illustration of the four features extracted for growth function fitting. A:** Time-domain responses at 60 dB SL. Shaded errorbars enclose three standard errors of the mean ( $n=151$  epochs). The peak to peak amplitude and the RMS value are extracted from a 50ms to 500ms post-stimulus window, **B:** Spectra of power change from baseline (arbitrary units, A.U.); black outline enclose the significant region of power change ( $z>3$ , permutation test with baseline at -300ms). Black dot indicates the time-frequency location of the peak power, and thus the extracted feature value. **C:** Spectra of phase locking value, significant region and peak location as with panel B, in arbitrary units (A.U.). **Bottom row: all subject responses to 60 dB SL stimuli. D:** All subjects' averaged time domain responses to 60 dB SL stimuli. Shaded errorbars enclose one standard deviation ( $n=20$  subjects). **E and F:** Power and phase locking spectra averaged across all subjects respectively. Black points signify peak location for each of the twenty subjects in the time-frequency domain and demonstrate that the spread of peak location across subjects is smaller for PLV.

**Figure 3: Each feature plotted against stimulus intensity over all subjects.** The filled circle and error bars represent the feature mean and one standard deviation respectively. Panels A to D indicate the peak-to-peak amplitude, root-mean-square, peak power and peak phase-locking-value features respectively.

**Figure 4: Growth function fitting procedure and threshold estimates using all epochs. A.** An example in subject 14 using the peak PLV to find the threshold estimate. **B.** Growth function fits (straight lines) and threshold estimates (circles) for all subjects. The black circles

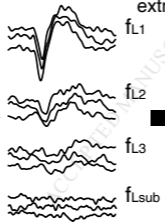
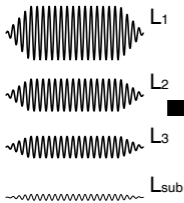
show the spread of threshold estimates across each subject. A point to the right of the vertical behavioral threshold line is a positive threshold estimate and indicates an over-estimation of the threshold, and vice versa for a negative estimate corresponding to an under-estimation. The horizontal bar indicates the mean and one standard deviation of threshold estimates across subjects. Four subjects were omitted from fits for peak power due to invalid fits (see text).

**Figure 5. Deterioration of threshold estimate accuracy when decreasing testing time. A.** Threshold estimate accuracy (quantified as standard deviation of threshold estimate) decreased with increasing PLV feature noise level, which is inversely related to acquisition time (Panel B). Each point is one threshold estimate and color indicates goodness of the growth function fit (Adjusted  $r^2$  value); red points indicate threshold estimate when using all epochs; black crosses are invalid fits. Black lines indicate threshold estimate mean and one standard deviation respectively. **B.** PLV feature noise decreases with increasing number of epochs and is largely independent of level. **C.** (Top) Scatterplot of threshold estimate standard deviation for each feature against their respective feature noise levels, an example of which is shown in panel A for peak PLV, before converting feature noise level to testing time in minutes. **D.** The percentage of invalid fits decreased when increasing test time.

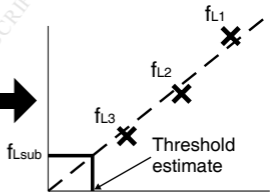
Sound stimulus

EEG activity

Growth function fit

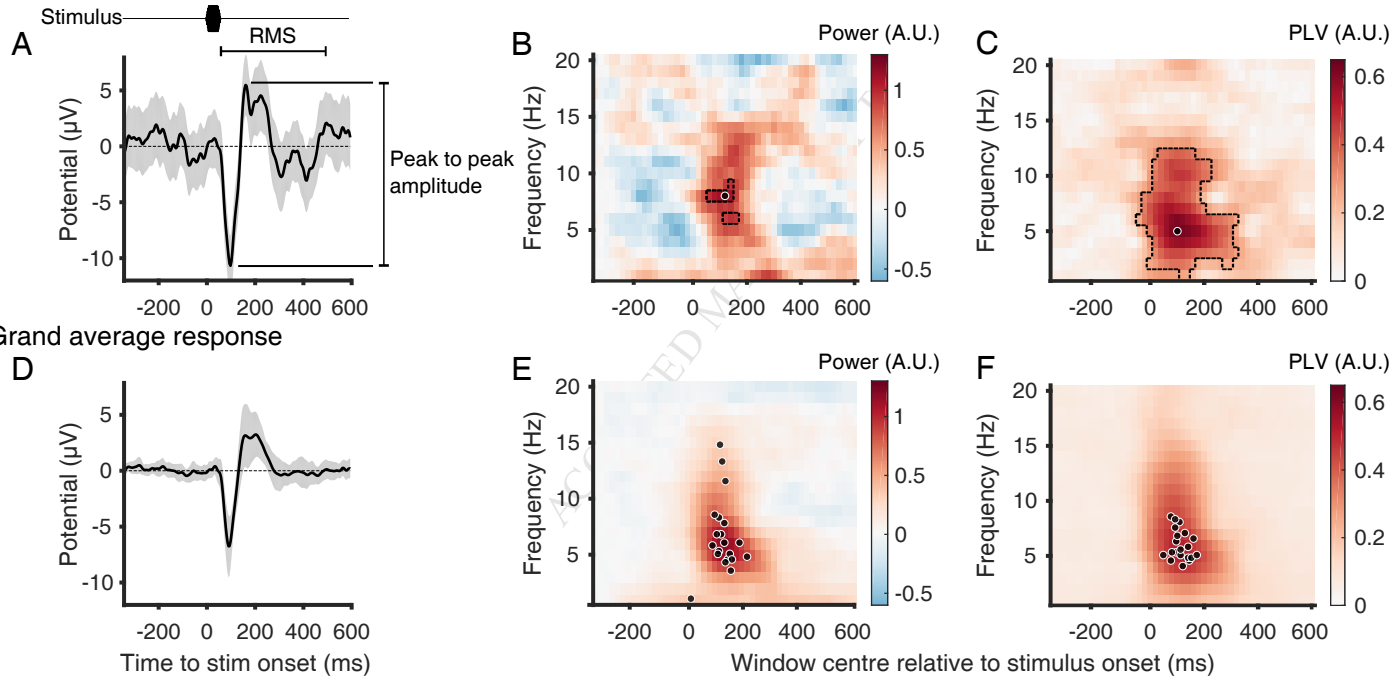


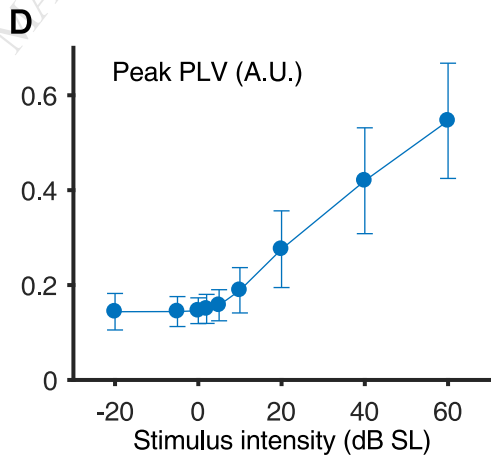
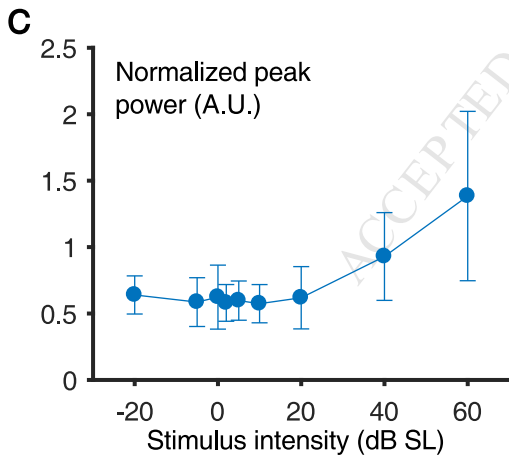
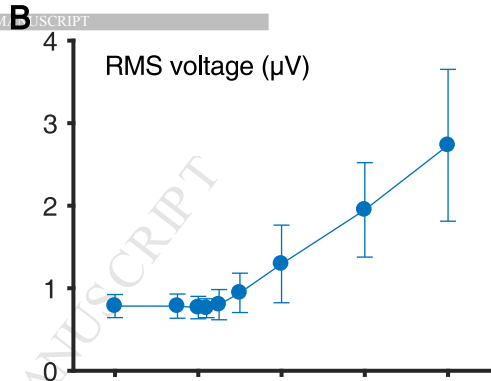
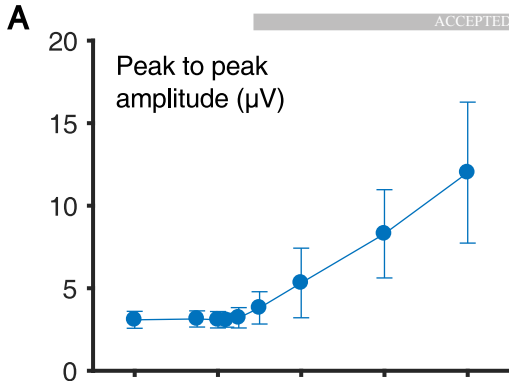
Feature extraction






Baseline activity gives  
baseline feature estimate

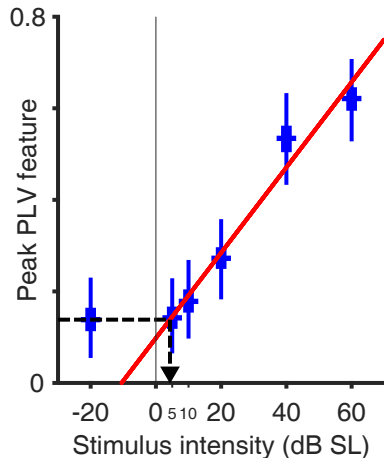
Sound stimulus intensity level





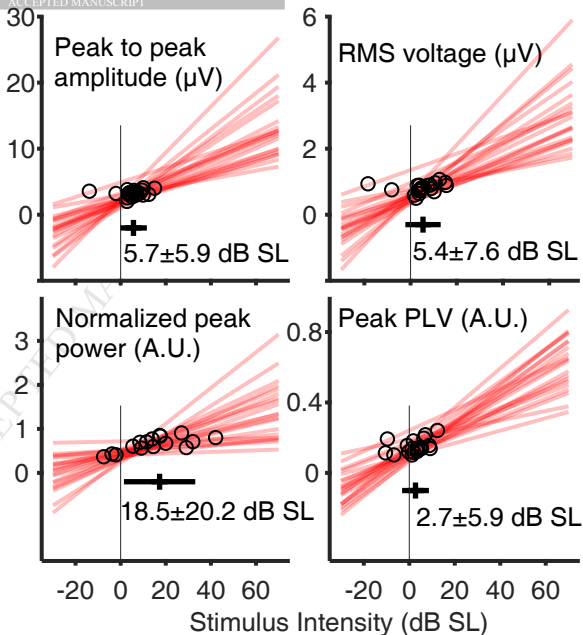
A

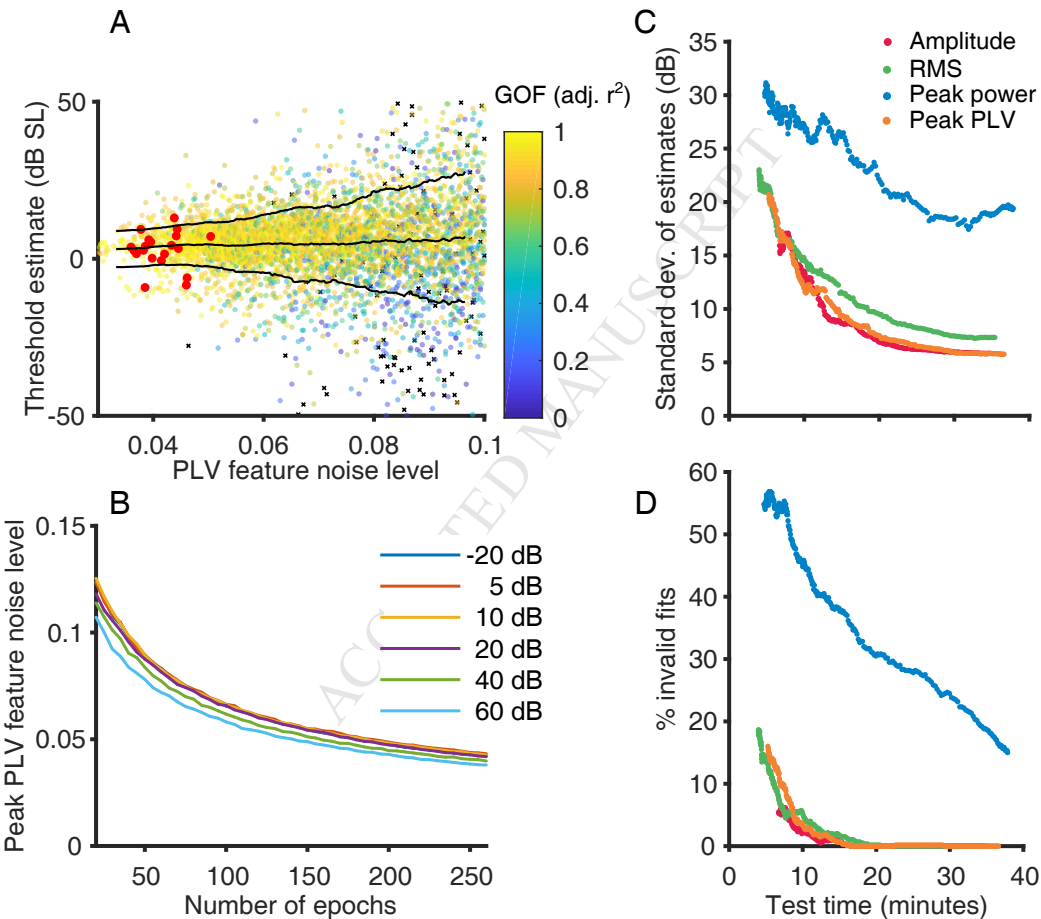
-  PLV bootstrap dist: Median, 25/75th & 1st/99th centiles
-  Linear curve fit
-  T estimate = 4.3 dB



B

ACCEPTED MANUSCRIPT





### Highlights

- Growth functions provide clinically viable objective estimates of hearing threshold
- Cortical response peak-to-peak amplitudes and phase locking values are suitable
- Best threshold estimates are consistently within 6 dB of behavioural threshold
- 13 minutes testing time for  $\pm 20$  dB accuracy with single electrode recording