

1 Comparing Phoneme and Word 2 Recognition Test Outcomes in Adult CI 3 users: Data Analysis from the AuDieT 4 Study

5 Enrico Migliorini^{1,2}, Jan-Willem A. Wasmann², Nikki Philpott², Bastiaan van Dijk³, Birgit
6 Philips¹, Wendy Huinck²

7 1) Cochlear Technology Centre Belgium, Mechelen, Belgium 2) Department of Otorhinolaryngology,
8 Donders Institute for Brain, Cognition and Behaviour, Radboud university medical center Nijmegen,
9 Nijmegen, The Netherlands 3) Cochlear Benelux NV, Mechelen, Belgium

10 Abstract

11 **Purpose:** Current clinical measures used in cochlear implantation (CI) provide a broader view of
12 speech recognition ability at word-level, often missing granular details contained at phoneme-level
13 that may be valuable for CI mapping. This study evaluates how outcomes of Phoneme Recognition in
14 Quiet tests (PRQ) differ from those of more commonly used word recognition tests (CVC) and
15 outlines how these tests may be useful for different purposes in clinical adult CI care.

16 **Methods:** As part of the AuDiET (Auditory Diagnostics and Error-based Treatment) study, 23 adult
17 postlingually deafened unilateral CI users underwent a battery of tests, including both PRQ and CVC
18 tests. Their results were compared at the phoneme level, including an evaluation of fitness and error
19 dispersion.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

20 **Results:** PRQ had a significantly lower accuracy and fitness than CVC. The error patterns also tended
21 to be less random and more systematic. Fitness correlated strongly and positively with accuracy,
22 while error dispersion negatively correlated with accuracy.

23 **Conclusion:** There are clear differences between PRQ and CVC outcomes in absolute accuracy and
24 error distribution. Comparing these tests might provide clinicians with more granular insights into
25 which areas/phonemes to target during mapping, to achieve optimal speech recognition.

26 Introduction

27 Background

28 The cochlear implant (CI) is a highly successful sensory neuro-prosthesis. Through electrical stimulation
29 of the auditory nerve, a CI can partially restore hearing in people with severe to profound sensorineural
30 hearing loss. Since its invention, it is estimated that more than a million hearing-impaired people
31 worldwide have received a CI (Zeng, 2022). Adult CI users score an average of 70%-80% speech
32 recognition in quiet when using the latest implants, sound processors, and coding strategies (Zeng,
33 2022).

34 While this is an impressive result, some issues related to the post-implantation journey of CI users still
35 need addressing. One of these is unexpectedly poor outcomes: cases in which a CI user achieves a
36 much lower level than expected (Pisoni et al., 2017). There are multiple possible explanations for this
37 variability: many of the factors which contribute to hearing, such as neural health in the spiral ganglion,
38 are hard to measure precisely (Walia et al., 2023), the impact of various factors on performance seems
39 to be variable (Blamey et al., 2013), and clinicians tend to overestimate the post-implantation
40 performance of CI recipients based on pre-implantation data (Philpott, Philips, Donders, et al., 2023).

41 Clinical research tends to focus on good or exceptional outcomes (Moberly et al., 2016), resulting in an
42 under-representation of people with poorer outcomes. There is a general lack of understanding on

43 how to address the issue of poor performance. This uncertainty affects the otherwise excellent cost-
44 effectiveness of CIs (Bond et al., 2009), and uncertainty about the potential benefits may dissuade
45 candidates from undergoing the implantation.

46 The issue of addressing CI users who show poor speech understanding is compounded by the lack of
47 patient-specific, standardised treatment guidelines. Whether considering fitting (i.e., adjusting the
48 parameters governing the electrical stimulation) or training (i.e., rehabilitative exercises aimed at
49 improving speech recognition), clinical practices vary significantly across different clinics (Wathour et
50 al., 2021). In clinical fitting procedures for the Cochlear Nucleus system, the two most important
51 settings which are individually determined are T-levels (just audible electrical stimulation level) and C-
52 levels (most comfortable electrical stimulation level). T-levels are determined using a threshold-
53 seeking method, while C-levels are set based on a CI user's loudness rating. Most parameters other
54 than T- and C-level are set to "default" (Vaerenberg et al., 2014). Surveys such as the aforementioned
55 Vaerenberg et al., 2014 show significant variability in CI programming practices among clinicians and CI
56 centres. The lack of standardised, evidence-based guidelines can lead different clinicians to adopt
57 different approaches, although steps have been taken towards developing the Living Guidelines:
58 consensus-based guidelines aimed at encouraging good clinical practices in clinics across the world.

59 While objective measures show a correlation with subjective threshold and comfort levels (de Vos et
60 al., 2018), the correlation is, in general, weak; this means that the objective measures cannot be used
61 to accurately predict T- and C-levels. However, it has been shown that creating MAPs based on either
62 electrically evoked Compound Action Potential (eCAP) or electrically evoked Stapedius Reflex
63 Threshold (eSRT) can lead to equivalent performance compared to behavioural subjective fitting
64 (Craddock et al., 2003).

65 In recent years, several studies have been conducted aiming to develop individualised fitting and
66 training interventions that take into account each participant's unique challenges. The investigation of

67 both fitting and training is important, as they represent different and complementary parts of the post-
68 implantation journey of a CI user. Most of the studies attempting to define a new fitting paradigm focus
69 on using objective, non-language-related measures to determine fitting levels: such tests may include
70 electrode discrimination, pitch ranking or spectrotemporal sensitivity tests (Grasmeder et al., 2019;
71 Van Opstal & Noordanus, 2023; Warren & Atcherson, 2023). Such “bottom-up” approaches stem from
72 the assumption that fittings that provide well-audible stimuli across the full spectrum of covered
73 frequencies would result in the best overall performance. Other approaches focus on imaging
74 techniques, aiming to fit on the base of anatomical features, electrode positioning, and distance from
75 the modiolus in order to leverage the tonotopic organization of the cochlea to prevent frequency drift
76 due to misalignment of the electrode with the spiral ganglion cells (Jiam et al., 2019; Kurz et al., 2023)
77 .

78 It is known that CI users react to different cues than people with typical hearing thresholds (Moberly et
79 al., 2014), which implies that delivering electrical stimulation with the aim of reproducing the neural
80 activity of a normal ear as accurately as possible (i.e., making it as faithful as possible to the actual
81 sound) may be a suboptimal approach. Enhancing the cues that CI users rely on, instead, may lead to
82 better speech understanding. However, there is a very limited number of studies concerning how
83 analysing speech recognition outcomes may provide insight into CI fitting (Holmes et al., 2012;
84 Wathour et al., 2023), several of which centre around the Fitting to Outcome eXperts (FOX; Otoconsult
85 NV, Antwerp, Belgium) system, an Artificial Intelligence (AI)-based fitting tool. The generally positive
86 outcome of papers relating to speech-based interventions suggests that there is merit to this approach:
87 however, the details of the interventions’ implementations have not been published in their entirety
88 due to both their commercial nature and, in the case of FOX, the intrinsically hard-to-understand
89 nature of AI. Therefore, the question of how clinicians may draft interventions based on their patients’
90 speech recognition outcomes remains partially unanswered.

91 As training interventions are complementary to fitting ones, it is worth mentioning how the concept of
92 individualised auditory training has also been investigated: Magits et al. (Magits et al., 2023) reported
93 similar effectiveness for personalised and non-personalised training programs, and in their literature
94 review, Philpott et al. (Philpott, Philips, Tromp, et al., 2023) found variable effectiveness for both types
95 of programs, although the personalisation of the training tended to refer to adaptive difficulty more
96 than to a focus on the individual difficulties of each CI user. In this case, as well as investigating
97 personalized fitting, we believed it worthwhile to evaluate whether personalized training based on the
98 speech recognition issues of each CI user may be impactful.

99 The Auditory Diagnostics and Error-based Treatment (AuDiET) study was conceived as an investigation
100 of the feasibility and effectiveness of individualised fitting and training interventions based on
101 phoneme-level information on the participants' errors. Its goal was to provide details on how CI users
102 experiencing different challenges in sound recognition, respond to interventions aimed at addressing
103 those challenges. Should an investigation of errors in speech audiometry provide valuable resources
104 upon which fitting and training interventions may be based, this could lead the way towards a new
105 approach to post-implantation clinical care. It may provide opportunities where users may test their
106 hearing performance on their own between visits (apps for remote testing with self-administered
107 procedures are already available, with varying degrees of functionality and validity (Wasmann et al.,
108 2024)) and detailed information on their errors could be relayed to their clinicians for use in their
109 follow-up.

110 In this paper we investigate whether there are significant differences in errors and error patterns when
111 comparing the results of phoneme tests and word tests (as described in the Methods section). Next,
112 differences between these tests are analysed to find out which test is better suited for designing
113 interventions. This will be done by analysing their correlation scores, highlighting differences in
114 accuracy scores, and dispersion both within participants and between the two different tests. Potential
115 explanations for the presented results will be considered in the Discussion section.

116 Methods

117 Study design

118 In the AuDiET study, each participant undergoes five clinical visits. During Visit 1, baseline data is
119 collected; during Visit 2, a fitting intervention is administered to the participant; during Visit 3 (set 2
120 weeks after Visit 2) the effects of the fitting intervention are evaluated, and the participant is given a
121 personalised training programme; during Visit 4 (set 4 weeks after Visit 3) the training intervention is
122 evaluated, and the participant stops training; finally, during Visit 5 (set 4 weeks after Visit 4) the
123 retention of any effects is evaluated.

124 The AuDiET study itself is structured as a pre-post comparison; the analysis presented in this paper,
125 however, is limited to the pre-intervention data collected from the study population. The study
126 population is comprised of 23 (27 were recruited, 4 of which dropped out of the study or were excluded
127 due to unforeseen technical issues) native Dutch-speaking adult CI users with a post-lingual onset of
128 hearing loss and unilaterally implanted with a Cochlear® Nucleus™ implant model. The details of the
129 population are highlighted in Table 1. All participants had at least one year of CI experience.
130 Participants with abnormally formed cochleae, severe pre-implantation ossification, severe cognitive
131 disorders, intense facial nerve stimulation, unaddressed tip fold-over or more than 4 malfunctioning
132 electrodes were ineligible for the study.

133 *Table 1: Information on the study population. 'S' stands for 'Subject', 'M' for 'Male', 'F' for 'Female'*

Subject No.	Age at inclusion (years range)	Gender	CI experience (y)	Implant type	Etiology
S01	61-65	M	5	CI522	Idiopathic
S02	71-75	M	7	CI512	Family history, otosclerosis

S03 (excluded)	-	-	-		
S04	66-70	F	4	CI522	Idiopathic
S05	66-70	F	4	CI422	Family history, idiopathic
S06	71-75	M	4	CI522	Idiopathic
S07	81-85	M	3	CI532	Otitis media
S08	86-90	M	5	CI522	Idiopathic
S09	76-80	F	6	CI522	Skull trauma
S10 (dropout)	-	-	-	-	-
S11	81-85	M	5	CI512	Family history, schwannoma
S12	76-80	M	8	CI422	Idiopathic
S13	66-70	M	5	CI522	Auditory neuropathy
S14	36-40	F	8	CI422	Congenital
S15	66-70	M	3	CI512	Sudden deafness, idiopathic
S16	76-80	M	12	CI512	Unknown
S17 (excluded)	-	-	-	-	-
S18	66-70	F	7	CI532	Family history
S19 (excluded)	-	-	-	-	-
S20	66-70	F	10	CI422	Family history, idiopathic
S21	56-60	F	12	CI24RE	Meningitis
S22	51-55	F	18	CI24RE	Vestibular schwannoma

S23	71-75	M	7	CI512	Skull trauma
S24	66-70	F	5	CI532	Congenital, rubella
S25	71-75	F	7	CI522	Radiotherapy for vestibular schwannoma
S26	66-70	F	4	CI422	DFNA-9
S27	76-80	M	7	CI522	Sudden deafness, idiopathic

134

135 The study was submitted to the ethical committee, where it was approved and assessed as not falling
136 under the jurisdiction of the Medical Research Involving Human Subjects Act (WMO) (NL-number:
137 NL80521.091.22). Recruitment and testing took place at Radboud university medical center between
138 2022 and 2023.

139 During Visit 1, each participant underwent a test battery aimed at collecting a detailed dataset of their
140 hearing capabilities. This battery included aided Pure Tone Audiometry (PTA), Spectrotemporal
141 Sensitivity Assessment (SSA, (Van Opstal & Noordanus, 2023)), Phoneme Recognition in Quiet (PRQ,
142 detailed below), Consonant-Vowel-Consonant (CVC, (Bosman & Smoorenburg, 1995)), and Digits
143 Triplet Test (Smits et al., 2013). All tests except Pure Tone Audiometry were streamed via Direct Audio
144 Cable (De Graaff et al., 2016) at a level of 65 input-related dBA to a Nucleus™ 6 test processor loaded
145 with the participant's most used MAP; Pure Tone Audiometry was instead performed in free-field using
146 the modified Hughson-Westlake staircase procedure and ensuring the blocking of the contralateral ear
147 in the case of residual hearing. The computer used for running all tests except PTA was a Lenovo
148 Thinkpad T440 (Lenovo, Hong Kong, Hong Kong), connected to a RME Fireface UC external sound card
149 (RME, Germany) to ensure consistent audio levels. Calibration was performed by directly reading the
150 DSP input levels with a tool provided by Cochlear Ltd. and comparing streamed levels vs levels as
151 acquired with the processor placed on a mannequin in a calibrated free field in a sound room. PTA was

152 performed according to clinical routine on a calibrated clinical audiometer in a sound booth or quiet
153 consultation room.

154 The data collected from the SSA and the DTT is not analysed and presented in this paper, as the goal
155 was to investigate differences between phoneme and word tests; the results of those two tests were
156 found to not be sufficiently granular for the purposes of investigating phoneme recognition. In the PRQ
157 test, participants listened to triphones of the form /hVt/ or /aCa/, where V represents vowels or
158 diphthongs (ɑ, a, au, ε, e, ei, ø, ɪ, i, ɔ, u, o, ʏ, œy, y) and C represents consonants (b, d, f, ɣ, h, j, k, l, m, n,
159 p, r, s, t, v, w, z) in the Dutch language. The participants were instructed to indicate what triphone they
160 heard from the closed set of all possible options (vowels and consonants tested separately). The full set
161 was presented in random order, 8 times for each consonant and 6 times for each vowel. The test
162 software was developed specifically for this study using Python 3.

163 In the CVC test, the participants heard 15 lists of 12 meaningful Dutch CVC words each, taken from the
164 Nederlandse Vereniging Audiologie (NVA) word list (Bosman & Smoorenburg, 1992), and were
165 required to type what was heard as a response. This was then automatically converted into a triplet of
166 phonemes in order to investigate errors on a phonemic level rather than judging words only as ‘correct’
167 or ‘incorrect’. The software for the CVC test was developed by Cochlear Ltd. and was previously
168 validated in a clinical study (de Graaff et al., 2018). For further information on the tests itself, refer to
169 the protocol as registered on <https://clinicaltrials.gov/study/NCT05307952>.

170 Both tests produced data points in the same format: arrays of stimulus-response pairings where the
171 stimulus is the phoneme being presented, and the response is the phoneme the participant reported
172 hearing. These arrays made up the primary outcome of the visit. In order to extract human-readable
173 information from these data points, several data transformation techniques were applied. These
174 include the calculation of accuracy (defined as the percentage of correctly identified phonemes),
175 fitness (also describable as weighted accuracy, described in detail below), and error dispersion (defined
176 in information theory as “the effective number of error classes per stimulus token” (Van Son, 1995)).

177 The measure of fitness was defined as $1 - d(stim, resp)$ where 'stim' is the presented phoneme and
178 'resp' is the response. The distance function d is defined as the perceptual distance between the
179 stimulus and the response, in such a way that if the stimulus and the response share some phonetic
180 features, they are marked as being closer than ones that differ completely. The features of phonemes
181 are those defined by the International Phonetic Association: voicing, place and manner for consonants,
182 openness, place and rounding for vowels (International Phonetic Association, 1999). For instance, the
183 distance between /p/ and /t/ can be set as 0.33 as they are both unvoiced plosives which only differ in
184 place, while the distance between /p/ and /z/ is the maximum of 1, as they differ in voicing, place, and
185 manner. In simple terms, fitness is lower when phonemes that are very different from each other are
186 confused. The details of the implementation can be found in the supplemental materials. It is worth
187 noting that, by definition, fitness dominates accuracy, i.e., for any given speech test, the fitness will
188 always be higher than the accuracy. This is because every error is marked as a zero when calculating
189 accuracy, while partial scores are possible when calculating fitness.

190 The data analysis aimed to compare the PRQ and CVC data, looking at correlations and differences
191 between phoneme and word tests. This was done first by calculating the correlation coefficients
192 between accuracy, fitness, and error dispersion in PRQ and CVC for each participant. The goal of this
193 was to highlight how accuracy, fitness, and error dispersion scores can be useful for describing
194 individual error patterns.

195 In order to assess within-visit test-retest reliability, random sampling was performed, splitting the data
196 for each visit in half and checking for a significant change in accuracy. The test was repeated ten times
197 and the results were averaged to reduce the chance of randomly selecting a split and returning a
198 spurious result.

199 Next, significant differences between the scores' distributions were investigated using the Wilcoxon
200 signed rank test for paired samples. Finally, the correlation between accuracy in vowels and
201 consonants was calculated for both PRQ and CVC. Results across the test are considered significant if

202 their p-values are lower than 0.05 after Bonferroni-Holm correction (performed using (Gaetano, 2013)
203). All reported p-values have been adjusted in this way.

204 Results

205 Overview and normality test

206 The means, medians, and standard deviations of the computed measures can be found in Table 2.
207 Using the Shapiro-Wilk test for normality, we found that the assumption of normality did not hold for
208 fitness in vowels, either for PRQ or CVC tests, and for the error dispersion of consonants in CVC tests.
209 For this reason, the non-parametric Wilcoxon test for paired samples was used instead of a parametric
210 one.

211 *Table 2: Aggregate measures and results of the Shapiro-Wilk test for each computed variable. Values below the 0.05*
212 *significance threshold are marked with an asterisk (*)*

Variable	Mean	Median	Standard deviation	P-value
Accuracy of PRQ vowels	69.13	72.06	16.49	0.09
Fitness of PRQ vowels	91.75	93.43	5.43	0.02*
Error dispersion of PRQ vowels	0.68	0.58	0.35	0.20
Accuracy of CVC vowels	82.49	86.78	13.45	0.07
Fitness of CVC vowels	95.84	96.80	3.90	<0.01*
Error dispersion of CVC vowels	1.09	1.10	0.70	0.19
Accuracy of PRQ consonants	66.62	69.12	19.45	0.33
Fitness of PRQ consonants	82.43	84.56	12.18	0.15
Error dispersion of PRQ consonants	0.96	0.88	0.69	0.02*
Accuracy of CVC consonants	78.56	76.21	13.51	0.06
Fitness of CVC consonants	89.09	88.85	7.12	0.11
Error dispersion of CVC consonants	1.57	1.71	0.76	0.68
Accuracy of PRQ vowels and consonants	67.59	71.00	16.81	0.35
Fitness of PRQ vowels and consonants	85.61	87.95	9.30	0.08

Error dispersion of PRQ vowels and consonants	0.84	0.76	0.51	0.08
Accuracy of CVC vowels and consonants	79.74	79.84	13.34	0.07
Fitness of CVC vowels and consonants	91.13	91.26	6.01	0.11
Error dispersion of CVC vowels and consonants	1.40	1.56	0.71	0.57

213

214 Correlation analysis – PRQ versus CVC scores

215 Error: Reference source not found shows the results of the CVC and PRQ tests for each subject. PRQ
 216 and CVC scores correlate strongly in all of accuracy (Pearson’s r : 0.90; p -value < 0.001), fitness
 217 (Pearson’s r : 0.87; p -value < 0.001) and error dispersion (Pearson’s r : 0.82; p -value < 0.001). These
 218 correlations remain present also when considering only vowels or only consonants.

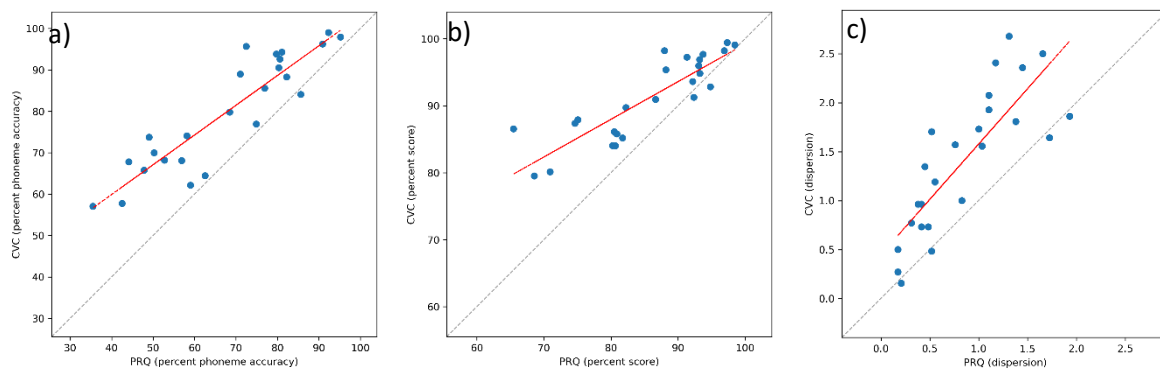


Figure 1: Distributions of PRQ and CVC scores for accuracy (a), fitness (b) and error dispersion (c). The red line indicates a linear regression through the data, the dotted grey line is the diagonal line where both scores are equal.

219 CVC tests results have a significantly higher accuracy (Wilcoxon’s p -value < 0.001), fitness (Wilcoxon’s
 220 p -value < 0.001) and error dispersion (Wilcoxon’s p -value < 0.001) than PRQ test results. The average
 221 accuracy, fitness, and error dispersion of CVC tests were 80%, 91%, and 1.4, respectively; those of PRQ
 222 tests were 68%, 86%, and 0.84.

223 In order to ensure that this difference in error dispersion was not linked to the number of phonemes
 224 presented in CVC being higher than that of PRQ, random sampling was used to repeat the test 10 times
 225 using a randomly selected number of CVC phonemes, equal in number and in vowels/consonants

226 percentage to the PRQ ones. The mean p -value of these ten randomly sampled tests was 0.044, which
227 is higher (due to the much reduced sample size) but still significant.

228 The random sampling aimed at investigating test-retest reliability reported no significant change in the
229 distribution of accuracy between any of the subsamples (the average p -value for a Wilcoxon test being
230 0.54).

231 Error: Reference source not found shows the data split between consonant and vowel scores.
232 Significant correlations between the accuracies of the consonants vs vowels were found both in PRQ
233 (Pearson's r : 0.57; p -value: 0.002) and CVC (Pearson's r : 0.95; p -value < 0.001).

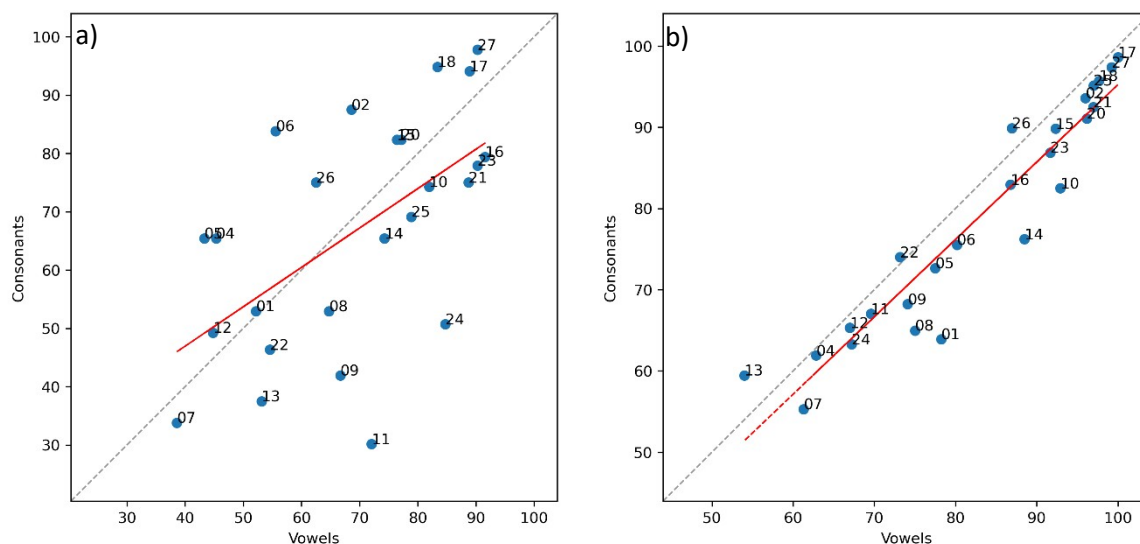


Figure 2: Distribution of accuracy in vowels and consonants for PRQ (a) and CVC (b) tests. The red line indicates a linear regression through the data, the dotted grey line is the diagonal line where both scores are equal. The numbers are Subject IDs.

234 Correlation analysis – Fitness and Error Dispersion

235 Figure 3 shows that fitness correlated strongly and positively with accuracy (Pearson's r : 0.97; p -value <
236 0.001).

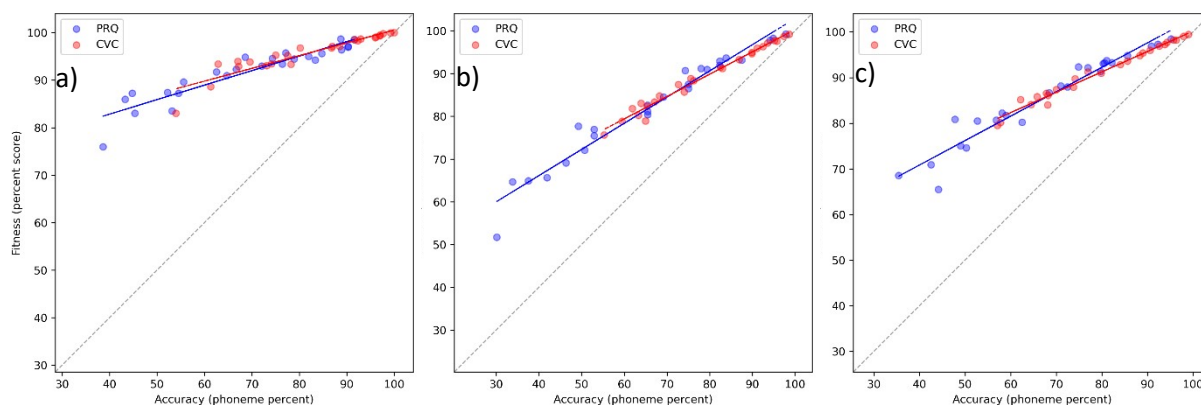


Figure 3: Distribution of fitness and accuracy for vowels (a), consonants (b) and overall (c). In these graphs the PRQ and CVC data is split by colour. The blue and red lines represent the linear regressions of PRQ and CVC data respectively. The dotted grey line is the diagonal line where both scores are equal.

237 Figure 4 shows that error dispersion correlates negatively with accuracy (Pearson's r : -0.61; p -value <
238 0.001).

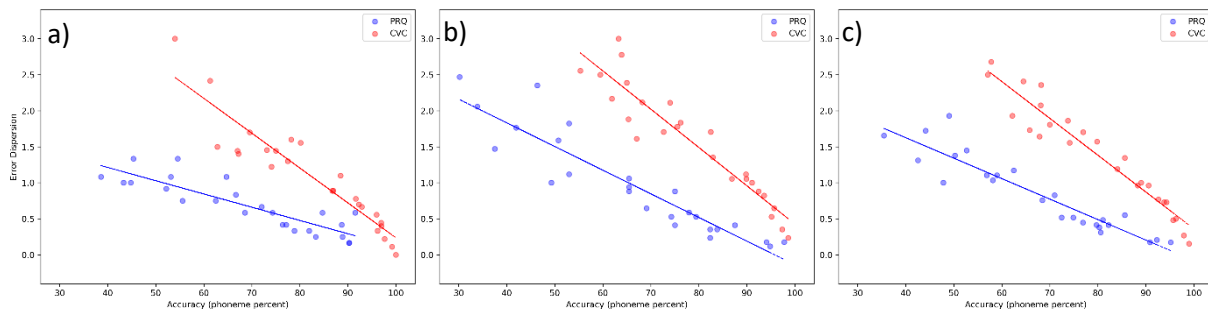


Figure 4: Distribution of error dispersion and accuracy for vowels (a), consonants (b) and overall (c). In these graphs the PRQ and CVC data is split by colour. The blue and red lines represent the linear regressions of PRQ and CVC data respectively. The dotted grey line is the diagonal line where both scores are equal.

239 Power analysis

240 Using the G*Power software (Faul et al., 2007) to calculate the achieved power, it was found to be
241 higher than 0.95 for each computer parameter, with the exception of the correlation between vowels
242 and consonants for PRQ, where the achieved power was 0.91. This is usually considered to be a very
243 high statistical power, resulting in a very low likelihood of a Type 1 error.

244 Discussion

245 In this paper we have shown how, within the study participants, Phoneme Recognition in Quiet test
246 scores differ from Consonant-Vowel-Consonants ones, having on average lower scores and more
247 consistent error patterns (as shown by the lower error dispersion scores).

248 We have also shown that different measures, such as fitness and error dispersion, can provide
249 additional information on a CI user's speech recognition; we believe that there is a strong case to be
250 made for the use of multiple scores (such as the aforementioned fitness and error dispersion) in clinical
251 practice, delivering more granular insights to clinicians on their patients' speech understanding. The
252 strong correlation of fitness with accuracy is expected as the former dominates the latter and both
253 measures are influenced by the number of errors made. However, in Figure 3, it is possible to see
254 participants whose results are placed at a certain distance from the regression line, and for these
255 comparing fitness and accuracy may be insightful. If the fitness of a participant's PRQ test is close to the

256 accuracy, it means that the participant is making large mistakes, confounding very different phonemes.
257 A deeper analysis of their errors might then provide insight into their issues and, according to our
258 hypothesis, make it so that an intervention could be designed to provide them with more easily
259 recognisable cues. Conversely, a larger difference between fitness and accuracy means that the
260 participant is confusing phonemes that are similar to each other, and they might be able to overcome
261 this issue with practice if the differences between the cues of the confused phonemes are too small to
262 intervene upon reliably. Furthermore, smaller errors (i.e., confusions between more similar
263 phonemes) should be easier for a subject to correct when making use of contextual cues, such as in a
264 meaningful sentence.

265 The negative correlation of error dispersion with accuracy can be explained by the fact that the fewer
266 errors a CI user is making, the lower the probability will be that those errors distribute into a high
267 number of categories. In a similar way to fitness, therefore, error dispersion is a measure that should
268 not be interpreted by itself, but in relation to accuracy.

269 The three scores taken together can provide clinicians with more detailed information on their
270 patients' individual issues. A participant who, for example, has difficulties distinguishing similar
271 phonemes such as /i/ and /e/ might be characterized by a fitness score significantly higher than
272 accuracy (since the phonemes being mistaken are similar) and a low error dispersion score (since the
273 errors are consistently made between those phonemes). Conversely, a participant who experiences
274 difficulties with recognizing phonemic cues would show a higher error dispersion score; those issues
275 might be addressed with training exercises aimed at practicing their phoneme discrimination and
276 recognition capabilities (Philpott, Philips, Tromp, et al., 2023).

277 Considering the mechanics of speech comprehension can help interpret the differences between these
278 tests. According to Erber's hierarchy (Erber, 1977), there are four steps involved in listening: Detection,
279 Discrimination, Recognition and Comprehension. By using meaningful words, CVC tests engage all skills
280 involved in each of these steps, including cognition. Instead, PRQ tests, which use meaningless

281 triphones, only involve the first three steps. CVC tests, by making use of cognition and top-down
282 processes and being influenced by co-articulation, may mask certain auditory errors and cause others
283 to appear. Subjects are told to expect meaningful words; therefore, they will correct perceptual errors,
284 increasing the accuracy of the test. Sometimes, however, they might introduce non-perceptual errors
285 by reporting a completely wrong, but meaningful word, e.g. (using English words for the sake of non-
286 Dutch-speaking readers), mishearing the word 'TOP' as 'TOG' and reporting hearing 'DOG' instead.

287 The different correlations of words with consonants in PRQ and CVC could also be interpreted in a
288 similar way. Considering that the CVC words are meaningful, correctly identifying the consonants in a
289 word provides additional cues for identifying the vowels and vice versa. Conversely, in PRQ the
290 participant can only rely on auditory cues. For instance, this could mean that a participant who
291 perceives spectral components well but has issues with temporal ones might perform well in vowel
292 recognition and more poorly in consonant recognition, which features temporal components more
293 prominently.

294 These results would suggest that PRQ is an effective way to evaluate how a CI user experiences speech
295 at a phonemic level, limiting the influence of co-articulation as well as cognition and comprehension
296 skills. Arguably, reducing the effect of cognition might benefit clinicians aiming to adjust CI users'
297 fitting, as the test can help identify bottom-up errors at a perceptual level. These low-level errors are
298 the ones which may be more effectively addressed by adjusting the CI user's fitting. In contrast, CVC
299 tests would be useful for replicating more closely speech in everyday life.

300 Another topic that may be further investigated is the implementation of phoneme tests at home. It
301 would be little effort to implement the PRQ test on a mobile application and let CI users perform it at
302 home; this would let clinicians have an overview of their patients' speech recognition issues without
303 devoting time to administering the test in a clinical environment. Research has shown that at-home
304 tests have the potential to be as reliable as those run in a clinic (van Wieringen et al., 2021; Wasmann
305 et al., 2024).

306 Finally, when considering how to assist a CI user optimally, the authors would recommend using a
307 combination of accuracy, fitness, and error dispersion for both PRQ and CVC tests. This approach aims
308 to paint a clearer picture of their individual difficulties in quiet. These scores might be further
309 integrated by data-savvy clinicians with Confusion Matrices to pinpoint the phonemes that each
310 participant has the most difficulties with.

311 [Limitations of the study](#)

312 While PRQ appears to be a reliable test to evaluate CI users' ability to identify phonemes without the
313 additional variability introduced by co-articulation and cognition, it is worthwhile to take some time to
314 discuss its limits. First, as the current study only included phonemes in quiet, results cannot be
315 translated to speech recognition in noise. A follow-up to this study investigating whether these results
316 hold for tests in noise is needed. Judging from previous studies (Goldsworthy et al., 2013), we can
317 expect both PRQ and CVC scores to deteriorate with the introduction of noise; however, there would
318 be merit in investigating which of the scores is most affected, and whether certain subsets of
319 phonemes are more impacted by noise.

320 Second, the responses being presented as a multiple-choice test might introduce a form of McGurk
321 effect (McGurk & MacDonald, 1976), inducing participants to report hearing phonemes that they read
322 but did not hear. Similarly, they may develop a subconscious bias, repeatedly choosing one phoneme
323 (for instance, the first one of the top row) when uncertain about what they heard.

324 Third, the test population consisted entirely of postlingually deafened, experienced CI users. A similar
325 study using prelingually deafened or newly implanted CI users might show different results. The study
326 by Magits et al. (Magits et al., 2023) included newly implanted subjects and found no differences
327 between inexperienced and experienced users, so it may be worthwhile to check whether these results
328 hold true for phoneme tests.

329 Finally, test-retest reliability needs further investigation in the context of multiple visits, and potentially
330 in the context of self-administered tests over a long period.

331 Conclusion

332 This paper presented Phoneme Recognition in Quiet testing as a valid integration to Consonant-Vowel-
333 Consonant Speech Audiometry testing. We showcased how the two tests seem to measure different
334 steps in the Erber hierarchy, respectively Recognition and Comprehension, and suggested that a
335 framework based on the use of multiple measures (accuracy, fitness, and error dispersion) over both
336 kinds of tests might provide audiologists with deeper insights into their patients' unique and individual
337 difficulties with speech recognition.

338 Further papers on the experimental parts of the AuDiET study will follow to investigate whether this
339 data-driven, individualized approach to fitting and training can improve the post-implantation follow-
340 up.

341 Ethics Declarations

342 Employment

343 Enrico Migliorini, Bastiaan van Dijk and Birgit Philips are employed by Cochlear Limited.

344 Funding

345 Enrico Migliorini's PhD programme is MOSAICS; MOSAICS is a European Industrial Doctorate project
346 funded by the European Union's Horizon 2020 framework programme for research and innovation
347 under the Marie Skłodowska-Curie Grant Agreement No. 860718.

348 Ethics approval

349 The study was carried on in conformity with the 1964 Declaration of Helsinki, the EU GDPR and all
350 applicable guidelines for the Kingdom of the Netherlands and Radboud university medical center.
351 The protocol was approved by Radboud university medical center's METC.

352 Informed Consent

353 All participants signed an informed consent module concerning their tests and gave explicit permission
354 for their anonymous data to be shared.

355 Bibliography

- 356 Blamey, P., Artieres, F., Bařkent, D., Bergeron, F., Beynon, A., Burke, E., Dillier, N., Dowell, R., Fraysse,
357 B., Gallégo, S., Govaerts, P. J., Green, K., Huber, A. M., Kleine-Punte, A., Maat, B., Marx, M.,
358 Mawman, D., Mosnier, I., O'Connor, A. F., ... Lazard, D. S. (2013). Factors affecting auditory
359 performance of postlinguistically deaf adults using cochlear implants: an update with 2251
360 patients. *Audiology & Neuro-Otology*, *18*(1), 36–47. <https://doi.org/10.1159/000343189>
- 361 Bond, M., Mealing, S., Anderson, R., Elston, J., Weiner, G., Taylor, R. S., Hoyle, M., Liu, Z., Price, A., &
362 Stein, K. (2009). The effectiveness and cost-effectiveness of cochlear implants for severe to
363 profound deafness in children and adults: a systematic review and economic model. *Health*
364 *Technology Assessment (Winchester, England)*, *13*(44). <https://doi.org/10.3310/HTA13440>
- 365 Bosman, A. J., & Smoorenburg, G. F. (1992). *Woordenlijst voor spraakaudiometrie*. Nederlandse
366 Vereniging Voor Audiologie, Utrecht. [https://scholar.google.com/scholar_lookup?](https://scholar.google.com/scholar_lookup?title=Woordeenlijst+Voor+Spraakaudiometrie&author=A.J.+Bosman&author=G.F.+Smoorenburg&publication_year=1992&)
367 [title=Woordeenlijst+Voor+Spraakaudiometrie&author=A.J.+Bosman&author=G.F.](https://scholar.google.com/scholar_lookup?title=Woordeenlijst+Voor+Spraakaudiometrie&author=A.J.+Bosman&author=G.F.+Smoorenburg&publication_year=1992&)
368 [+Smoorenburg&publication_year=1992&](https://scholar.google.com/scholar_lookup?title=Woordeenlijst+Voor+Spraakaudiometrie&author=A.J.+Bosman&author=G.F.+Smoorenburg&publication_year=1992&)
- 369 Bosman, A. J., & Smoorenburg, G. F. (1995). Intelligibility of Dutch CVC syllables and sentences for
370 listeners with normal hearing and with three types of hearing impairment. *Audiology : Official*
371 *Organ of the International Society of Audiology*, *34*(5), 260–284.
372 <https://doi.org/10.3109/00206099509071918>
- 373 Craddock, L., Cooper, H., van de Heyning, P., Vermeire, K., Davies, M., Patel, J., Cullington, H., Ricaud,
374 R., Brunelli, T., Knight, M., Plant, K., Cafarelli Dees, D., & Murray, B. (2003). Comparison
375 between NRT-based MAPs and behaviourally measured MAPs at different stimulation rates--a
376 multicentre investigation. *Cochlear Implants International*, *4*(4), 161–170.
377 <https://doi.org/10.1179/CIM.2003.4.4.161>
- 378 de Graaff, F., Huysmans, E., Merkus, P., Theo Goverts, S., & Smits, C. (2018). Assessment of speech
379 recognition abilities in quiet and in noise: a comparison between self-administered home
380 testing and testing in the clinic for adult cochlear implant users. *International Journal of*
381 *Audiology*, *57*(11), 872–880. <https://doi.org/10.1080/14992027.2018.1506168>
- 382 De Graaff, F., Huysmans, E., Qazi, O. U. R., Vanpoucke, F. J., Merkus, P., Goverts, S. T., & Smits, C.
383 (2016). The Development of Remote Speech Recognition Tests for Adult Cochlear Implant
384 Users: The Effect of Presentation Mode of the Noise and a Reliable Method to Deliver Sound in

- 385 Home Environments. *Audiology & Neuro-Otology*, 21 Suppl 1(1), 48–54.
386 <https://doi.org/10.1159/000448355>
- 387 de Vos, J. J., Biesheuvel, J. D., Briaire, J. J., Boot, P. S., van Gendt, M. J., Dekkers, O. M., Fiocco, M., &
388 Frijns, J. H. M. (2018). Use of Electrically Evoked Compound Action Potentials for Cochlear
389 Implant Fitting: A Systematic Review. *Ear and Hearing*, 39(3), 401–411.
390 <https://doi.org/10.1097/AUD.0000000000000495>
- 391 Erber, N. (1977). Evaluating speech-perception ability in hearing impaired children. In *Childhood*
392 *Deafness* (pp. 173–181). Grune & Stratton.
- 393 Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G*Power 3: a flexible statistical power
394 analysis program for the social, behavioral, and biomedical sciences. *Behavior Research*
395 *Methods*, 39(2), 175–191. <https://doi.org/10.3758/BF03193146>
- 396 Gaetano, J. (2013). *Holm-Bonferroni sequential correction: an Excel Calculator*.
- 397 Goldsworthy, R. L., Delhorne, L. A., Braid, L. D., & Reed, C. M. (2013). Psychoacoustic and Phoneme
398 Identification Measures in Cochlear-Implant and Normal-Hearing Listeners. *Trends in*
399 *Amplification*, 17(1), 27. <https://doi.org/10.1177/1084713813477244>
- 400 Grasmeder, M. L., Verschuur, C. A., van Besouw, R. M., Wheatley, A. M. H., & Newman, T. A. (2019).
401 Measurement of pitch perception as a function of cochlear implant electrode and its effect on
402 speech perception with different frequency allocations. *International Journal of Audiology*,
403 58(3), 158–166. <https://doi.org/10.1080/14992027.2018.1516048>
- 404 Holmes, A. E., Shrivastav, R., Krause, L., Siburt, H. W., & Schwartz, E. (2012). Speech based
405 optimization of cochlear implants. *International Journal of Audiology*, 51(11), 806–816.
406 <https://doi.org/10.3109/14992027.2012.705899>
- 407 International Phonetic Association. (1999). *Handbook of the International Phonetic Association: A*
408 *guide to the use of the International Phonetic Alphabet*. Cambridge University Press.
- 409 Jiam, N. T., Gilbert, M., Cooke, D., Jiradejvong, P., Barrett, K., Caldwell, M., & Limb, C. J. (2019).
410 Association Between Flat-Panel Computed Tomographic Imaging-Guided Place-Pitch Mapping
411 and Speech and Pitch Perception in Cochlear Implant Users. *JAMA Otolaryngology-- Head &*
412 *Neck Surgery*, 145(2), 109–116. <https://doi.org/10.1001/JAMAOTO.2018.3096>
- 413 Kurz, A., Herrmann, D., Hagen, R., & Rak, K. (2023). Using Anatomy-Based Fitting to Reduce
414 Frequency-to-Place Mismatch in Experienced Bilateral Cochlear Implant Users: A Promising
415 Concept. *Journal of Personalized Medicine*, 13(7). <https://doi.org/10.3390/JPM13071109>
- 416 Magits, S., Boon, E., De Meyere, L., Dierckx, A., Vermaete, E., Francart, T., Verhaert, N., Wouters, J., &
417 Van Wieringen, A. (2023). Comparing the Outcomes of a Personalized Versus Nonpersonalized
418 Home-Based Auditory Training Program for Cochlear Implant Users. *Ear and Hearing*, 44(3),
419 477–493. <https://doi.org/10.1097/AUD.0000000000001295>
- 420 McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature* 1976 264:5588,
421 264(5588), 746–748. <https://doi.org/10.1038/264746a0>
- 422 Moberly, A. C., Bates, C., Harris, M. S., & Pisoni, D. B. (2016). The Enigma of Poor Performance by
423 Adults with Cochlear Implants. *Otology & Neurotology : Official Publication of the American*
424 *Otological Society, American Neurotology Society [and] European Academy of Otology and*
425 *Neurotology*, 37(10), 1522. <https://doi.org/10.1097/MAO.0000000000001211>
- 426 Moberly, A. C., Lowenstein, J. H., Tarr, E., Caldwell-Tarr, A., Welling, D. B., Shahin, A. J., & Nittrouera,
427 S. (2014). Do adults with cochlear implants rely on different acoustic cues for phoneme

- 428 perception than adults with normal hearing? *Journal of Speech, Language, and Hearing*
429 *Research : JSLHR*, 57(2), 566. https://doi.org/10.1044/2014_JSLHR-H-12-0323
- 430 Philpott, N., Philips, B., Donders, R., Mylanus, E., & Huinck, W. (2023). Variability in clinicians'
431 prediction accuracy for outcomes of adult cochlear implant users. *International Journal of*
432 *Audiology*. <https://doi.org/10.1080/14992027.2023.2256973>
- 433 Philpott, N., Philips, B., Tromp, K., Kramer, S., Mylanus, E., & Huinck, W. (2023). Phoneme Training for
434 Adult Cochlear Implant Users: A Review of the Literature and Study Protocol. *Journal of Speech,*
435 *Language, and Hearing Research : JSLHR*, 66(12), 5071–5086.
436 https://doi.org/10.1044/2023_JSLHR-23-00335
- 437 Pisoni, D. B., Kronenberger, W. G., Harris, M. S., & Moberly, A. C. (2017). Three challenges for future
438 research on cochlear implants. *World Journal of Otorhinolaryngology - Head and Neck Surgery*,
439 3(4), 240. <https://doi.org/10.1016/J.WJORL.2017.12.010>
- 440 Smits, C., Theo Goverts, S., & Festen, J. M. (2013). The digits-in-noise test: assessing auditory speech
441 recognition abilities in noise. *The Journal of the Acoustical Society of America*, 133(3), 1693–
442 1706. <https://doi.org/10.1121/1.4789933>
- 443 Vaerenberg, B., Smits, C., De Ceulaer, G., Zir, E., Harman, S., Jaspers, N., Tam, Y., Dillon, M., Wesarg,
444 T., Martin-Bonniot, D., Gärtner, L., Cozma, S., Kosaner, J., Prentiss, S., Sasidharan, P., Briaire, J.
445 J., Bradley, J., Debruyne, J., Hollow, R., ... Govaerts, P. J. (2014). Cochlear implant programming:
446 a global survey on the state of the art. *TheScientificWorldJournal*, 2014.
447 <https://doi.org/10.1155/2014/501738>
- 448 Van Opstal, A. J., & Noordanus, E. (2023). Towards personalized and optimized fitting of cochlear
449 implants. *Frontiers in Neuroscience*, 17. <https://doi.org/10.3389/FNINS.2023.1183126>
- 450 Van Son, R. (1995). A method to quantify the error distribution in confusion matrices. *EUROSPEECH*.
- 451 van Wieringen, A., Magits, S., Francart, T., & Wouters, J. (2021). Home-Based Speech Perception
452 Monitoring for Clinical Use With Cochlear Implant Users. *Frontiers in Neuroscience*, 15.
453 <https://doi.org/10.3389/FNINS.2021.773427>
- 454 Walia, A., Shew, M. A., Lefler, S. M., Ortmann, A. J., Durakovic, N., Wick, C. C., Herzog, J. A., &
455 Buchman, C. A. (2023). Factors Affecting Performance in Adults With Cochlear Implants: A Role
456 for Cognition and Residual Cochlear Function. *Otology & Neurotology : Official Publication of*
457 *the American Otological Society, American Neurotology Society [and] European Academy of*
458 *Otology and Neurotology*, 44(10), 988–996. <https://doi.org/10.1097/MAO.0000000000004015>
- 459 Warren, S. E., & Atcherson, S. R. (2023). Evaluation of a clinical method for selective electrode
460 deactivation in cochlear implant programming. *Frontiers in Human Neuroscience*, 17.
461 <https://doi.org/10.3389/FNHUM.2023.1157673>
- 462 Wasmann, J. W. A., Huinck, W. J., & Lanting, C. P. (2024). Remote Cochlear Implant Assessments:
463 Validity and Stability in Self-Administered Smartphone-Based Testing. *Ear and Hearing*, 45(1),
464 239–249. <https://doi.org/10.1097/AUD.0000000000001422>
- 465 Wathour, J., Govaerts, P. J., & Deggouj, N. (2021). Variability of fitting parameters across cochlear
466 implant centres. *European Archives of Oto-Rhino-Laryngology : Official Journal of the European*
467 *Federation of Oto-Rhino-Laryngological Societies (EUFOS) : Affiliated with the German Society*
468 *for Oto-Rhino-Laryngology - Head and Neck Surgery*, 278(12), 4671–4679.
469 <https://doi.org/10.1007/S00405-020-06572-W>
- 470 Wathour, J., Govaerts, P. J., Lacroix, E., & Naïma, D. (2023). Effect of a CI Programming Fitting Tool
471 with Artificial Intelligence in Experienced Cochlear Implant Patients. *Otology and Neurotology*,

472 44(3), 209–215. <https://doi.org/10.1097/MAO.0000000000003810>

473 Zeng, F.-G. (2022). Celebrating the one millionth cochlear implant. *JASA Express Letters*, 2(7).
474 <https://doi.org/10.1121/10.0012825>

475

476 Supplemental Material

477 Fitness calculation

478 Let t be a test made of n presentations of phonemes. Since t can be described as n couples of
479 phonemes $\{p, q\}$ where p is the presented phoneme and q is the given answer. The IPA features of p
480 and q (voicing, place, and manner for consonants; rounding, place, and openness for consonants) are
481 then compared. If p and q do not overlap in any of the features, the $\{p, q\}$ couple is scored 0. If they
482 overlap in one, the couple is scored $1/3$; if they overlap in two, it is scored $2/3$, and if they overlap
483 completely then $p = q$, and the couple is scored 1. The average score over the n couples in the test t is
484 the fitness of t .