

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

Who is Right?: A word identification in noise test for young children using minimal pair distracters

Samuel Evans¹ & Stuart Rosen²

¹ Department of Psychology, University of Westminster, 115 New Cavendish St, London, W1W 6UW, UK

² Department of Speech, Hearing & Phonetic Sciences, UCL, 2 Wakefield St., London, WC1N 1PF

31 **Abstract (200 words)**

32

33 Many children have difficulties understanding speech. Reliable tools are needed to
34 identify them in order to provide appropriate interventions. At present, there are
35 relatively few assessments that test for subtle impairments in speech perception that
36 have normative data from UK children. Here we present a new test, which evaluates
37 children's ability to identify target words in background noise by choosing between
38 minimal pair alternatives that differ by a single articulatory phonetic feature. This
39 new test of single word perception is (1) tailored to testing young children, (2) has
40 minimal memory demands, (3) adapts to the child's ability and (4) does not require
41 reading or verbal output. Although designed for young children, it is also readily
42 applicable in adults. Here, we show that speech in noise abilities in this particular
43 task develop rapidly through childhood until they reach maturity at around ten years
44 of age. We make this test freely available, with normative data for listeners aged 4-
45 25 years old, and hope that it will be useful to researchers and clinicians in the
46 assessment of speech perceptual abilities in children with hearing impairments,
47 Developmental Language Disorder (DLD), Dyslexia and Auditory Processing
48 Disorder (APD).

49

50 **Key words: Speech perception, children, noise, audiological testing**

51

52

53 **Introduction**

54 Children with speech, language and hearing disorders are at a greater risk of
55 poorer literacy (Anthony & Francis, 2005), psycho-social development (Kilpatrick et
56 al., 2019) and long term prospects (Bryan et al., 2007). Deficits in speech
57 perception, in addition to being a defining feature of hearing impairment and Auditory
58 Processing Disorder (APD) (Moore et al., 2013), are associated with a number of
59 developmental disorders, most notably dyslexia (Noordenbos & Serniclaes, 2015)
60 and Developmental Language Disorder (DLD) (Ferguson et al., 2011). Developing
61 robust methods to identify individuals with speech perceptual deficits is a first step
62 towards better characterising and treating these disorders. At present, there are few
63 tests of speech perception that assess subtle impairments in speech perception and
64 that have appropriate normative data from UK children. Here, we make freely
65 available such a test, which we envisage will be useful to researchers and clinicians
66 in evaluating the perceptual abilities of young children.

67 Many children find understanding spoken language difficult. In some children,
68 such as those with hearing impairments, these difficulties are obvious and affect
69 perception in both ideal and adverse listening situations. Pure tone thresholds,
70 although important, provide limited information on functional listening abilities
71 (Houtgast & Festen, 2008) and tests of speech perception in noise provide arguably
72 a more valid assessment of day-to-day listening in children (Leibold et al., 2019).
73 Children with developmental language disorders can exhibit subtle speech
74 perceptual deficits, which often only become apparent when listening is made more
75 challenging (Moore et al., 2013; Ziegler et al., 2005). Indeed, deficits are not always
76 readily apparent and are sometimes only found in a minority of individuals, or not at
77 all (Messaoud-Galusi et al., 2011). This may reflect a lack of sensitivity of available

78 tests, an absence of a true speech perception deficit or significant heterogeneity in
79 the individuals assigned to these groups - only further research will help to uncover
80 which of these explanations is correct. This task is made more difficult by the high
81 co-morbidity between developmental reading, language and auditory processing
82 disorders (Bishop et al., 2016; Moore et al., 2013) and the paucity of tools for
83 assessing speech perception in children. A wider range of speech perception tests
84 are required to better characterise the speech perceptual abilities of hearing
85 impaired children and to further our understanding of developmental language
86 disorders.

87 Successful speech perception requires the integration of multiple co-varying
88 acoustic features (Kluender & Alexander, 2010; Lisker, 1977). In natural speech, the
89 multiplicity of available features helps to ensure that perception remains relatively
90 robust to acoustic variation and degradation of the speech signal. Speech sounds
91 that differ on the basis of fewer contrastive features are more highly confusable
92 (Miller & Nicely, 1955). Children with language impairments tend to perform more
93 poorly on tasks in which speech tokens differ minimally from one another
94 (Zoubrinetzky et al., 2016) and less so in tasks involving natural speech tokens that
95 differ on the basis of multiple acoustic cues (Coady et al., 2005). Speech perception
96 tasks can also be made more challenging by manipulating extrinsic factors, such as
97 the presence of competing noise. Competing sounds generate overlapping patterns
98 of excitation in the auditory periphery that obscure salient acoustic cues, a
99 phenomena referred to as energetic/modulation masking (Brungart, 2001; Stone et
100 al., 2011). Additional, informational masking effects, those not explained by
101 energetic and modulation masking, are thought to arise at more central, cognitive
102 levels of processing (Shinn-Cunningham, 2008). This form of masking is most often

103 associated with competing speech and is attributable in part to the difficulty of
104 separating out and attending to the correct speech stream. Difficulties attending to
105 speech in both competing speech (Dole et al., 2012) and non-speech maskers
106 (Ziegler et al., 2005, 2009) have been associated with developmental language
107 disorders.

108 Unfortunately, there are relatively few tests of speech perception in noise
109 designed specifically for children with normative data from the UK. Tests designed
110 for children need to be made engaging and use appropriate linguistic materials. It is
111 important that tests have normative data from the country in which they are used.
112 Normative data from other English speaking countries is unlikely to be appropriate
113 for use in the UK and can sometimes overestimate the prevalence of perceptual
114 deficits (Dawes & Bishop, 2007). Tests such as the SCAN-C (Keith, 2000) have
115 been adapted for use with British children (Dawes & Bishop, 2007). However,
116 SCAN-C is arguably not ideal for testing children with language impairments as it
117 requires them to repeat back heard words. Many children with language disorders
118 have difficulty planning and producing speech (Bishop et al., 2016) and so tests that
119 require a verbal response may underestimate their true abilities.

120 For the same reason, tests such as the FAAF that require children to read
121 words (Foster & Haggard, 1987) and those using sentences (e.g. LISN-S, Cameron
122 & Dillon, 2007) that place greater demands on auditory working memory and
123 syntactic processing, may not always be appropriate. Sentence material may be
124 particularly inappropriate given the evidence that sentence repetition in *quiet*
125 appears to be a good way to diagnose DLD (Conti-Ramsden et al., 2001). Tests that
126 use single, early acquired words and that require a non-verbal output response

127 provide a purer test of speech perceptual ability in those with speech and language
128 impairments.

129 There are relatively few existing UK tests of single word perception that have
130 a non-verbal output response. The Consonant Confusion Test (CCT) is suitable for
131 very young children and requires them to identify a target word from 4 alternatives
132 presented as pictures. However, in this test the alternatives differ by multiple
133 phonemes, e.g. “cow, owl, house, mouse”, hence the degree of phonemic
134 discrimination required in this task is relatively broad. The Chear Auditory
135 Perception Test (CAPT) is appropriate for slightly older children and includes
136 contrasts that require a finer level of discrimination. However, the normative data for
137 both these tests are derived from presenting the words at an artificially low volume,
138 used as a way of inducing variation in accuracy (Vickers et al., 2018). This is
139 arguably a less ecologically valid approach, compared to using competing noise to
140 bring accuracy ‘off ceiling’.

141 The McCormick Toy Test (Summerfield et al., 1994) combines phonemic
142 discrimination with concurrent noise presentation. However, the phonemic contrasts
143 between word alternatives are not always minimal (e.g., “man” vs. “lamb”). Vance et
144 al. (2009) includes fine grained phonemic discriminations, such that many of the
145 items differ on a single articulatory phonetic feature, with concurrent noise
146 presentation. However, the use of a fixed rather than an adaptive noise level does
147 not accommodate children performing at the extremes of accuracy. Indeed, this kind
148 of variation in performance is more likely in heterogeneous samples like those with
149 developmental language disorders.

150 Here, we present a new speech perception test, the Who is Right? (WiR?)
151 test and associated normative data for UK children and young adults. In this
152 computer administered task, the listeners identify a target spoken word from three
153 spoken alternative utterances that are presented against a competing noise.
154 Participants indicate their response non-verbally with a button press. To ensure
155 maximum sensitivity in identifying subtle impairments of speech processing, these
156 alternatives differ by a single articulatory phonetic feature, with the background noise
157 level adjusted adaptively dependent on their trial to trial performance accuracy.

158 **Methods & Materials**

159 **Test construction**

160 The WiR consists of 42 trials, all of a similar form. On each trial, the listener is
161 presented with a picture of a target word on a display screen and hears a male
162 speaker produce the name of the target in quiet (see Figure 1). Below the picture of
163 the target are three cartoon faces which then take turns to speak three utterances.
164 These utterances are produced by a different, female voice, to prevent participants
165 using an echoic memory trace to complete the task. The voices are presented
166 against a background of speech-spectrum-shaped noise (see details below). Two of
167 the utterances are non-word foils differing from the target in its initial consonant in a
168 single feature of voicing, place or manner (with the two foils always differing in the
169 contrast used). The other face produces the target. For example, when the target is
170 “bed”, the foils are “med” (differing in manner) and “ped” (differing in voicing). The
171 position of the target and two distracter foils are randomised from trial to trial. The
172 listener’s task is to identify the face that produced the correct target word by clicking
173 on that face using a mouse. A correct response results in the selected cartoon face

174 smiling, whereas an incorrect response results in the selected face frowning. Every
175 test began with a presentation of 14 familiarisation items followed by 28 test items,
176 with a random permutation of the items within each phase. All stimuli were
177 presented over headphones at a fixed comfortable level of about 65 dB SPL
178 (measured over the frequency range 100 Hz – 5 kHz).

179 Targets words were early-acquired monosyllabic words mainly of CVC
180 structure (two targets are CVs), that could be presented in an unambiguous pictorial
181 form and whose initial consonant could be altered by a single feature of voicing,
182 manner or place, to create two non-word foils (see Supplementary Materials for full
183 details). All items were early acquired words, and the test items had a mean age of
184 acquisition of 4.0 years, ranging from 2.9 to 5.6 (sd = 0.67), as measured by
185 Kuperman et al. (2012). For the test trials, the distracter foils comprised 14 manner
186 change items, 21 place change items and 21 voicing change items, distributed over
187 the 28 test trials (2 feature changes per target).

188 During the test, the signal-to-noise ratio (SNR) was varied adaptively using a
189 two-down/one-up adaptive rule tracking 71% correct (Levitt, 1971), which means that
190 the SNR increases after every error, and decreases after two consecutive correct
191 responses. The starting SNR was 20 dB, with a step-size of 7 dB which decreased
192 by 1 dB after every track reversal until it reached 3 dB, at which value it remained for
193 the rest of the test. The Speech Reception Threshold (SRT), here defined as the
194 SNR that led to about 71% correct responses, was calculated from the mean of the
195 track reversals in the test phase. Note that lower values indicate better performance,
196 as this indicates that the listener can tolerate poorer SNRs for the desired accuracy.

197 The stimulus triplets differed greatly in inherent intelligibility, as would be
198 expected by their variety of acoustic, phonetic and psycholinguistic properties. This
199 is highly undesirable in adaptive testing because it leads to greater variability in the
200 adaptive track. Extensive prior testing on dozens of school-age children allowed the
201 determination of the psychometric functions (relating proportion correct to SNR) for
202 each individual triplet. SRTs were then derived from these functions (through logistic
203 regression) allowing the calculation of a correction factor (the deviation for each
204 triplet from the mean SRT for all triplets) that was applied to the nominal SNR
205 desired during each test. In this way, performance should be similar for all triplets at
206 the same nominal SNR, which leads to more stable estimates of the SRTs.

207 The foil triplets were presented against a background of speech-spectrum-
208 shaped shaped noise, synthesised to approximate the long-term average speech
209 spectrum for combined male and female voices as estimated from the study of Byrne
210 et al., (1994). This consisted of a low-frequency portion rolling off below 120 Hz at
211 17.5 dB/octave, and a high-frequency portion rolling off at 7.2 dB/octave above 420
212 Hz, with a constant spectrum portion in-between. The noise started 450 ms before
213 the utterance triplet and finished 250 ms after, running continuously through the
214 three utterances, with 50 ms rise and fall times. The test, including all materials, and
215 analyses presented in this article are available here:

216 <https://github.com/drstuartrosen/WholsRight>.

217 **Participants**

218 Ethical approval was granted by the UCL Research Ethics Committee.
219 Informed written consent was received from all participants, and their parents, for
220 those aged less than 16 years. None of the children or adults tested had any known

221 speech, hearing or language impairments and they were all native British English
222 speakers.

223 The children and young adults were tested in primary and secondary schools
224 in four separate rounds of testing – referred to as SC (n = 30), GY (n =17), RL (n =
225 54) and HR (n =17) – and were combined in the analysis. In all instances, testing
226 took place in a quiet room either within school, home or in a quiet, distraction free
227 public space, in Southern England. One round of testing (GY) arose from control
228 data from typically developing children as part of a broader study of developmental
229 language disorder (Baird et al., 2011; Loucas et al., 2016).

230 There were 118 participants who completed the test and for whom there was
231 complete demographic information (mean age = 12.7 years, ranging from 4.0 to
232 25.1, with a s.d. = 4.63). Gender was well balanced with 53 males and 65 females
233 (55%).

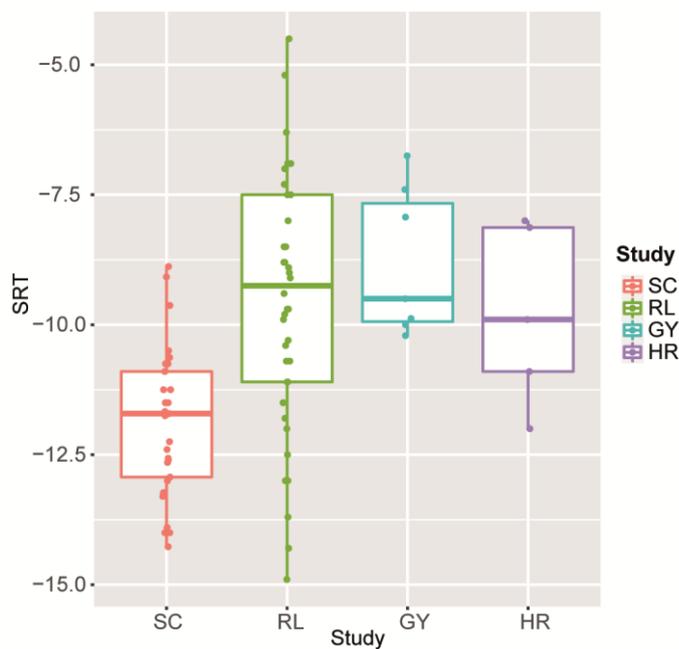
234 **Results**

235 The mean over the reversals in the test phase of the adaptive track was used
236 to estimate a *Speech Reception Threshold* (SRT) for each participant. A plot of the
237 obtained data against age shows a strong developmental trend of improving SRTs
238 up to about age 10 with a significant levelling off after that point. This also
239 suggested that the SRTs from the SC group that mainly included older participants
240 were on average better than the other groups for participants of a similar age. There
241 was one clear outlying low SRT ($z = -3$) which was deleted.

242 On the basis of the evidence that SRTs did not improve after age 11,
243 boxplots were made of the SRTs from the 4 studies for all listeners greater than that

244 age (Figure 1). A one-way ANOVA with a follow-up Tukey post-hoc test confirmed
245 the observation that the mean SRTs were not the same across the 4 testing groups
246 ($f(3, 77) = 9.203, p = 2.825e-05$). The SRTs for SC were significantly different from
247 RL and GY (both adjusted $ps < 0.003$), but SC and HR were not significantly different
248 from each other ($p = 0.152$) even though the absolute difference in means was very
249 similar to the other two groups, which did differ.

250



251

252 **Figure 1: SRTs for children aged 11 years and above in 4 different testing**
253 **rounds, illustrating the lower SRTs in the SC study.**

254

255 It is not clear why SRTs were lower in this group and we assume that this
256 reflects random sampling error. As SC only had participants aged 11.6-16.5 years
257 (in secondary school), it seemed undesirable to leave the SRTs as they were,
258 because the overall effect on model fits would not be equal across the age range.
259 Therefore, all SRTs in the SC study were adjusted by the mean difference between
260 the SRTs in that study and the three other studies for children ≥ 11 years old only

261 (i.e., by 2.4 dB). A one-way ANOVA confirmed that there was no evidence for
262 differences across the groups after the adjustment ($f(3,77) = 0.301$, $p = 0.825$).

263 On the evidence that SRTs change up to about age 10, and then asymptote,
264 two different models were used to fit the data. One was a segmented, or broken stick
265 regression, in which the model consists of two straight lines which meet at a
266 breakpoint. In this fit, a model in which the upper line had a slope=0 after the
267 breakpoint (implying no change in SRTs after a particular age), was statistically
268 indistinguishable from a model with non-zero slope for the upper segment. The
269 estimated breakpoint was estimated at 9.9 years (95% CI = 8.6 – 11.2). Note that,
270 for completeness, the data were also analysed without the adjustment accounting for
271 the lower SRTs in the SC study and the findings were similar, with a breakpoint at
272 age 10.4 years.

273 The other model was an asymptotic regression model with the equation:

274

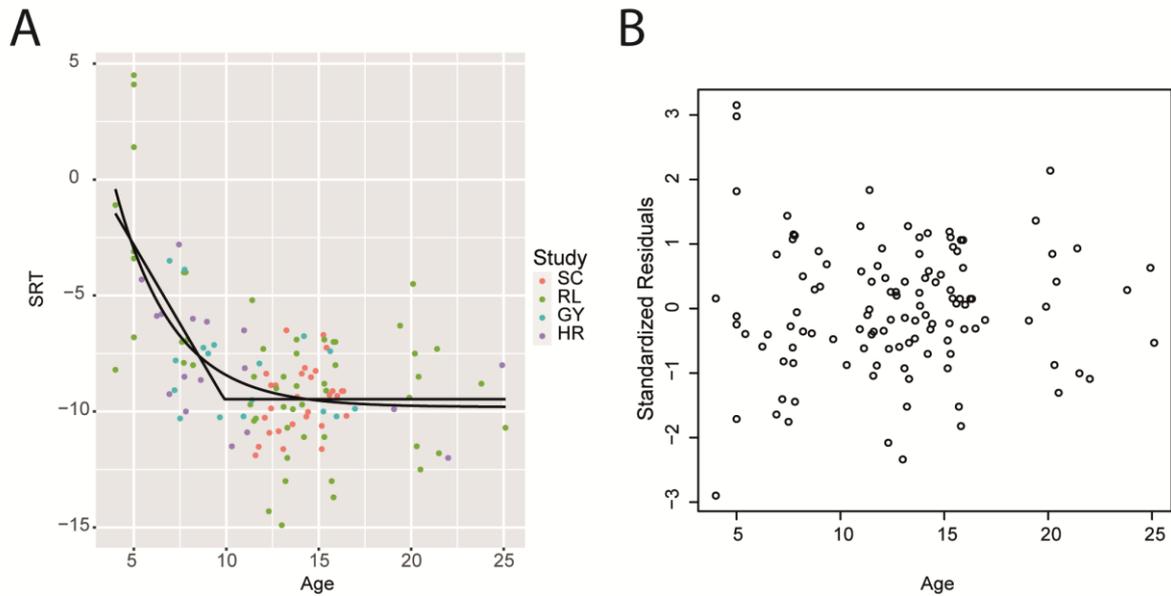
$$275 \text{SRT} = b_1 + b_2 * \exp(b_3 * \text{age})$$

276

277 where b_1 represents the asymptotic value (i.e., the lowest SRT reached through
278 development), as long as $b_3 < 0$, which was indeed the case; b_3 controls how fast
279 SRTs change over age, and b_2 scales the total range of this change. Note the
280 important interaction between b_2 and b_3 in determining the shape of the curve,
281 whereas b_1 is a simple additive term.

282

283



284

285 **Figure 2: Regression models of SRT with age. A) SRT regressed against age**
286 **with data sets indicated showing the broken stick and asymptotic regression**
287 **fits. B) Standardized residuals of the broken stick regression.**

288

289 The overall fits of the two models were very similar, as shown in Figure 2A,
290 with the broken stick model a slightly better fit with a residual standard error of 2.34
291 compared to 2.40 on 114 degrees of freedom (as both models have the same
292 number of estimated parameters). Visualisation of the standardised residuals
293 against age for the broken stick regression indicated that variability in measurement
294 of SRT was relatively constant across age after 5 years (Figure 2B).

295 As for many diagnostic tests, instead of expressing the outcome in a unit that
296 a test directly manipulates (here, SNR in dB), it is often more useful to calculate a z-
297 score, which reflects an individual's level of performance in comparison to their age-
298 matched peers. This is straightforward to do based on the broken stick regression.

299

300 First, a predicted SRT must be calculated based on the listener's age, where:

301 *If age ≤ 9.9, Predicted SRT = -1.36 x age + 3.98*

302 *If age > 9.9, Predicted SRT = - 9.5*

303

304 Then, a residual is calculated by subtracting the predicted SRT from the
305 actual SRT. This indicates by how many dB a listener is better or worse than an age-
306 matched peer, with negative numbers again indicating better performance. This is
307 then expressed as a z-score by dividing by an estimate of the standard deviation of
308 the residuals (2.33). From the z-score, a percentile can be calculated.

309 Suppose, for example, that a child aged 6 obtained an SRT of -0.7 dB. The
310 predicted SRT would be -4.2 dB from the equation above, which means this child is
311 3.5 dB worse than expected. Dividing through by 2.33 gives $z = 1.5$, which is to say,
312 1.5 standard deviations worse than typical 6 year olds. Only about 7% of children of
313 that age would be expected to have an SRT this poor or worse. SRT values in dB
314 and as z-scores relative to this normative data is provided as output for users of the
315 test.

316 **Discussion**

317 We have presented normative data from UK children on a test of word
318 identification in noise using minimal pair distracters. A broken stick regression
319 showed that perceptual abilities on this task continued to improve rapidly until the
320 age of around 10 years, before levelling out. We make this task and associated
321 normative data freely available and hope that this test will be of use to researchers
322 and clinicians in the assessment of speech perception abilities of children with

323 hearing, speech and language impairments. In the following sections, we discuss
324 future developments and limitations of the task.

325 Native language speech sound representations are relatively well developed
326 by 24 months of age but continue to be further refined well into later childhood (Kuhl,
327 2011). However, the point at which they achieve full maturity is still unknown.
328 Changes are observed until at least six years of age (Nittrouer & Studdert-Kennedy,
329 1987; Nittrouer, 2002) with some studies showing that maturation continues beyond
330 the early teens (Hazan & Barrett, 2000) and into the late teenage years (Davis et al.,
331 2019; McMurray et al., 2018) . In the WiR? test, performance rapidly improves until
332 around 10 years, before reaching a plateau. This break point is very similar to that
333 obtained in a similar open response word recognition task in speech-spectrum-noise
334 in a US sample (Corbin et al., 2016) and is broadly aligned with other studies
335 showing rapid development of speech in noise abilities up until the age of ten for
336 tasks involving competing energetic maskers (Hall et al., 2002; Leibold & Buss,
337 2013; Nishi et al., 2010; Wightman & Kistler, 2005). The earlier maturation on this
338 task, compared to the tasks described above in which maturation continues into the
339 late teenage years, may be attributed to important task differences. Our task
340 requires participants to discriminate between canonical articulations with perceptual
341 ambiguity arising from an extrinsic source, the presence of competing noise. By
342 contrast, categorical perception paradigms require participants to categorise
343 ambiguous sounds that are synthesised to be intermediate between canonical
344 articulations. This may require a finer level of phonetic discrimination or place
345 differing demands on decision making and executive function that gives rise to a
346 different developmental trajectory.

347 The early plateau in energetic masking abilities stands in contrast to the more
348 protracted development associated with informational masking. Indeed, children
349 perform worse with speech on speech masking compared to equivalent non-speech
350 maskers. Adult like performance on these tasks is not achieved until much later,
351 often beyond 13 years of age (Corbin et al., 2016; Hall et al., 2002; Leibold & Buss,
352 2013). Future developments of the WiR could include the incorporation of
353 competing speech maskers to provide sensitivity to the effects of information
354 masking, especially in older children. In addition, the facility to spatially separate
355 target and masking sounds, with associated normative data, may be a useful
356 inclusion to identify those with APD that have a specific difficulty with spatial sound
357 processing (Cameron & Dillon, 2007). Furthermore, given the different minimal pair
358 contrasts available in WiR?, it may be possible to collect normative data on specific
359 phonetic contrasts within the test. The ability to differentiate the contrasts that
360 children find most difficult may provide a perspective on the mechanisms that
361 underlie their speech perceptual weaknesses and allow better targeted interventions
362 for children with hearing impairments. However, it is likely that such tests would
363 require a fixed SNR, rather than an adaptive approach, with the SNR being fixed at a
364 level appropriate for the listener. In this way, it could be assured that listeners would
365 be not performing near floor or ceiling, but obtain intermediate levels of performance
366 which would allow meaningful comparisons across contrast types.

367 The task in its current form also has limitations. At present we do not have a
368 measure of re-test reliability or an understanding of how performance on the test
369 changes with repetitive testing. We hope that re-test reliability would be relatively
370 high given the efforts made to calibrate the task through the estimation of an SNR
371 correction factor for each item. Visualisation of the standardised residuals of our

372 normative data showed that they are relatively uniformly distributed with few outliers
373 suggesting that the SRT measure is relatively stable across age. We anticipate that
374 learning in the task would be minimal both within a single test session and across
375 multiple tests due to the relatively large number of test words and the fact that they
376 are not repeated. Future work addressing re-test reliability and learning effects will
377 help to clarify whether our intuitions are correct.

378 Another limitation is in the normative data that we have acquired. Our full
379 sample was 117 participants, a sample size roughly in keeping with or better than
380 similar tests (Spyridakou et al., 2020; Vance et al., 2009; Vickers et al., 2018).
381 However, as with most tools of this kind, it would benefit from a larger normative
382 sample, from a broader demographic, as factors like social economic status have
383 been shown to influence speech perceptual ability (Nittrouer, 1996). Indeed, our
384 data was collected from only a small number of settings and likely represents a
385 relatively homogenous demographic sample. In future, normative data from a wider
386 demographic including hard to reach populations is necessary, taking into account
387 the additional time and resources that this would entail (Bonevski et al., 2014). As
388 part of this widening inclusion, it would also be beneficial to consider stratifying by
389 UK region to account for differences in regional accent (Adank et al., 2009). We
390 would hope to address these limitations in the future and to allow others to do so, by
391 making this test freely available, so that others can extend upon our initial work.

392 **Acknowledgements**

393 We are very grateful to Gillian Baird and Vicky Slonims for providing the GY data
394 from their much larger study. More importantly, it was the desire expressed by their
395 team for an appropriate test of speech-in-noise to be used in a study of children with

396 language disorders that was the prime impetus for the development of WiR? The
397 rest of the data was extracted from the BSc theses of Sarah Cooper, Rebecca
398 Lancaster and Henrietta Louise Roe, to whom we are also grateful. Samuel Evans
399 was funded by a Vacation Scholarship from the charity Deafness Research UK,
400 which was merged with the RNID (Royal National Institute for Deaf People).

401

402 **Declaration of interests**

403 None

404

405 **References**

- 406 Adank, P., Evans, B. G., Stuart-Smith, J., & Scott, S. K. (2009). Comprehension of
407 Familiar and Unfamiliar Native Accents Under Adverse Listening Conditions.
408 *Journal of Experimental Psychology: Human Perception and Performance*,
409 *35*(2), 520–529. <https://doi.org/10.1037/a0013552>
- 410 Anthony, J., & Francis, D. (2005). Development of Phonological Awareness skill.
411 *Current Directions in Psychological Science*, *14*(5), 255–259.
412 http://www.speechpathology.com/articles/article_detail.asp?article_id=342
- 413 Baird, G., Slonims, V., Simonoff, E., & Dworzynski, K. (2011). Impairment in non-
414 word repetition: A marker for language impairment or reading impairment?
415 *Developmental Medicine and Child Neurology*, *53*(8), 711–716.
416 <https://doi.org/10.1111/j.1469-8749.2011.03936.x>
- 417 Bishop, D. V. M., Snowling, M. J., Thompson, P. A., Greenhalgh, T., Adams, C.,
418 Archibald, L., Baird, G., Bauer, A., Bellair, J., Boyle, C., Brownlie, E., Carter, G.,
419 Clark, B., Clegg, J., Cohen, N., Conti-Ramsden, G., Dockrell, J., Dunn, J.,
420 Ebbels, S., ... Whitehouse, A. (2016). CATALISE: A multinational and
421 multidisciplinary Delphi consensus study. Identifying language impairments in
422 children. *PLoS ONE*, *11*(7), 1–26. <https://doi.org/10.1371/journal.pone.0158753>
- 423 Bonevski, B., Randell, M., Paul, C., Chapman, K., Twyman, L., Bryant, J., Brozek, I.,
424 & Hughes, C. (2014). Reaching the hard-to-reach: A systematic review of
425 strategies for improving health and medical research with socially
426 disadvantaged groups. *BMC Medical Research Methodology*, *14*(1), 1–29.
427 <https://doi.org/10.1186/1471-2288-14-42>
- 428 Brungart, D. S. (2001). Informational and energetic masking effects in the perception
429 of two simultaneous talkers. *Journal of the Acoustical Society of America*,
430 *109*(3), 1101–1109. wos:000167369300023

- 431 Bryan, K., Freer, J., & Furlong, C. (2007). Language and communication difficulties
 432 in juvenile offenders. *International Journal of Language and Communication*
 433 *Disorders*, 42(5), 505–520. <https://doi.org/10.1080/13682820601053977>
- 434 Cameron, S., & Dillon, H. (2007). Development of the listening in spatialized noise-
 435 sentences test (LISN-S). *Ear and Hearing*, 28(2), 196–211.
 436 <https://doi.org/10.1097/AUD.0b013e318031267f>
- 437 Coady, J. A., Kluender, K. R., & Evans, J. L. (2005). Categorical perception of
 438 speech by children with specific language impairments. *Journal of Speech,*
 439 *Language, and Hearing Research*, 48(4), 944–959.
 440 [https://doi.org/10.1044/1092-4388\(2005/065\)](https://doi.org/10.1044/1092-4388(2005/065))
- 441 Conti-Ramsden, G., Botting, N., & Faragher, B. (2001). Psycholinguistic markers for
 442 specific language impairment (SLI). *Journal of Child Psychology and Psychiatry*
 443 *and Allied Disciplines*, 42(6), 741–748. <https://doi.org/10.1111/1469-7610.00770>
- 444 Corbin, N. E., Bonino, A. Y., Buss, E., & Leibold, L. J. (2016). Development of open-
 445 set word recognition in children: Speech-shaped noise and two-talker speech
 446 maskers. *Ear and Hearing*, 37(1), 55–63.
 447 <https://doi.org/10.1097/AUD.0000000000000201>
- 448 Davis, M., Evans, S., McCarthy, K., Evans, L., Giannakopoulou, A., & Taylor, J.
 449 (2019). Lexical learning shapes the development of speech perception until late
 450 adolescence. *PsyArXiv*, 1–52. <https://doi.org/10.31234/osf.io/ktsey>
- 451 Dawes, P., & Bishop, D. V. M. (2007). The SCAN-C in testing for auditory processing
 452 disorder in a sample of British children. In *International Journal of Audiology*
 453 (Vol. 46, Issue 12). <https://doi.org/10.1080/14992020701545906>
- 454 Dole, M., Hoen, M., & Meunier, F. (2012). Speech-in-noise perception deficit in
 455 adults with dyslexia: Effects of background type and listening configuration.
 456 *Neuropsychologia*, 50(7), 1543–1552.
 457 <http://www.sciencedirect.com/science/article/pii/S0028393212001200>
- 458 Ferguson, M. A., Hall, R. L., Riley, A., & Moore, D. R. (2011). Communication,
 459 listening, cognitive and speech perception skills in children with auditory
 460 processing disorder (APD) or specific language impairment (SLI). *Journal of*
 461 *Speech, Language, and Hearing Research*, 54(1), 211–227.
 462 [https://doi.org/10.1044/1092-4388\(2010/09-0167\)](https://doi.org/10.1044/1092-4388(2010/09-0167))
- 463 Foster, J. R., & Haggard, M. P. (1987). The Four Alternative Auditory Feature test
 464 (FAAF)-linguistic and psychometric properties of the material with normative
 465 data in noise. *British Journal of Audiology*, 21(3), 165–174.
 466 <https://doi.org/10.3109/03005368709076402>
- 467 Hall, J. W., Grose, J. H., Buss, E., & Dev, M. B. (2002). Spondee recognition in a
 468 two-talker masker and a speech-shaped noise masker in adults and children.
 469 *Ear and Hearing*, 23(2), 159–165. [https://doi.org/10.1097/00003446-200204000-](https://doi.org/10.1097/00003446-200204000-00008)
 470 [00008](https://doi.org/10.1097/00003446-200204000-00008)
- 471 Hazan, V., & Barrett, S. (2000). The development of phonemic categorization in
 472 children aged 6-12. *Journal of Phonetics*, 28(4), 377–396.

- 473 <https://doi.org/10.1006/jpho.2000.0121>
- 474 Houtgast, T., & Festen, J. M. (2008). On the auditory and cognitive functions that
475 may explain an individual's elevation of the speech reception threshold in noise.
476 *International Journal of Audiology*, 47(6), 287–295.
477 <https://doi.org/10.1080/14992020802127109>
- 478 Keith, R. W. (2000). Development and standardization of SCAN-C test for auditory
479 processing disorders in children. *Journal of the American Academy of*
480 *Audiology*, 11(8), 438–445.
- 481 Kilpatrick, T., Leitão, S., & Boyes, M. (2019). Mental health in adolescents with a
482 history of developmental language disorder: The moderating effect of bullying
483 victimisation. *Autism and Developmental Language Impairments*, 4.
484 <https://doi.org/10.1177/2396941519893313>
- 485 Kluender, K. R., & Alexander, J. M. (2010). Perception of Speech Sounds. In A. I.
486 Basbaum, A. Kaneko, G. M. Shephard, G. Westheimer, T. D. Albright, R. H.
487 Masland, P. Dallos, D. Oertel, S. Firestein, G. K. Beauchamp, C. Bushnell, J. H.
488 Kass, & E. Gardner (Eds.), *The Senses: A Comprehensive Reference*. Elsevier.
- 489 Kuhl, P. (2011). Brain Mechanisms in Early Language Acquisition. *Neuron*, 67(5),
490 713–727.
- 491 Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition
492 ratings for 30,000 English words. *Behavior Research Methods*, 44(4), 978–990.
493 <https://doi.org/10.3758/s13428-012-0210-4>
- 494 Leibold, L. J., & Buss, E. (2013). Children's identification of consonants in a speech-
495 shaped noise or a two-talker masker. *Journal of Speech, Language, and*
496 *Hearing Research*, 56(4), 1144–1155. [https://doi.org/10.1044/1092-](https://doi.org/10.1044/1092-4388(2012/12-0011))
497 [4388\(2012/12-0011\)](https://doi.org/10.1044/1092-4388(2012/12-0011))
- 498 Leibold, L. J., Buss, E., & Calandruccio, L. (2019). Too Young for the Cocktail Party?
499 *Acoustics Today*, 15(1), 37. <https://doi.org/10.1121/at.2019.15.1.39>
- 500 Lisker, L. (1977). Rapid Versus Rabid - Catalog of Acoustic Features That May Cue
501 Distinction. *Journal of the Acoustical Society of America*, 62, S77–S78.
502 [isi:A1977EA29000377](https://doi.org/10.1121/1.381111)
- 503 Loucas, T., Baird, G., Simonoff, E., & Slonims, V. (2016). Phonological processing in
504 children with specific language impairment with and without reading difficulties.
505 *International Journal of Language and Communication Disorders*, 51(5), 581–
506 588. <https://doi.org/10.1111/1460-6984.12225>
- 507 McMurray, B., Danelz, A., Rigler, H., & Seedorff, M. (2018). Speech categorization
508 develops slowly through adolescence. *Developmental Psychology*, 54(8), 1472–
509 1491. <https://doi.org/10.1037/dev0000542>
- 510 Messaoud-Galusi, S., Hazan, V., & Rosen, S. (2011). Investigating speech
511 perception in children with dyslexia: Is there evidence of a consistent deficit in
512 individuals? *Journal of Speech, Language, and Hearing Research*, 54(6), 1682–
513 1701. [https://doi.org/10.1044/1092-4388\(2011/09-0261\)](https://doi.org/10.1044/1092-4388(2011/09-0261))

- 514 Miller, G., & Nicely, P. (1955). An analysis of perceptual confusions among some
515 English consonant. *Journal of the Acoustical Society of America*, 27, 338–352.
516 <https://doi.org/10.1121/1.1907526>
- 517 Moore, D. R., Rosen, S., Bamiou, D. E., Campbell, N. G., Sirimanna, T., James
518 Bellis, T., Chermak, G., Weihing, J., Musiek, F., Dillon, H., & Cameron, S.
519 (2013). Evolving concepts of developmental auditory processing disorder (APD):
520 A British Society of Audiology APD Special Interest Group “white paper.”
521 *International Journal of Audiology*, 52(1), 3–13.
522 <https://doi.org/10.3109/14992027.2012.723143>
- 523 Nishi, K., Lewis, D. E., Hoover, B. M., Choi, S., & Stelmachowicz, P. G. (2010).
524 Children’s recognition of American English consonants in noise. *The Journal of*
525 *the Acoustical Society of America*, 127(5), 3177–3188.
526 <https://doi.org/10.1121/1.3377080>
- 527 Nittrouer, S., & Studdert-Kennedy, M. (1987). The role of coarticulatory effects in the
528 perception of fricatives by children and adults. *Journal of Speech and Hearing*
529 *Research*, 30(3), 319–329. <https://doi.org/10.1044/jshr.3003.319>
- 530 Nittrouer, Susan. (1996). The relation between speech perception and phonemic
531 awareness: Evidence from low-SES children and children with chronic OM.
532 *Journal of Speech, Language, and Hearing Research*, 39(5), 1059–1070.
533 <https://doi.org/10.1044/jshr.3905.1059>
- 534 Nittrouer, Susan. (2002). Learning to perceive speech: How fricative perception
535 changes, and how it stays the same. *The Journal of the Acoustical Society of*
536 *America*, 112(2), 711–719. <https://doi.org/10.1121/1.1496082>
- 537 Noordenbos, M. W., & Serniclaes, W. (2015). The Categorical Perception Deficit in
538 Dyslexia: A Meta-Analysis. *Scientific Studies of Reading*, 19(5), 340–359.
539 <https://doi.org/10.1080/10888438.2015.1052455>
- 540 Shinn-Cunningham, B. G. (2008). Object-based auditory and visual attention. *Trends*
541 *in Cognitive Sciences*, 12(5), 182–186. wos:000256750400006
- 542 Spyridakou, C., Rosen, S., Dritsakis, G., & Bamiou, D. E. (2020). Adult normative
543 data for the speech in babble (SiB) test. *International Journal of Audiology*,
544 59(1), 33–38. <https://doi.org/10.1080/14992027.2019.1638526>
- 545 Stone, M. a, Füllgrabe, C., Mackinnon, R. C., & Moore, B. C. J. (2011). The
546 importance for speech intelligibility of random fluctuations in “steady”
547 background noise. *The Journal of the Acoustical Society of America*, 130(5),
548 2874–2881. <https://doi.org/10.1121/1.3641371>
- 549 Summerfield, Q., Palmer, A. R., Foster, J. R., Marshall, D. H., & Twomey, T. (1994).
550 Clinical evaluation and test-retest reliability of the IHR-McCormick Automated
551 Toy Discrimination Test. *British Journal of Audiology*, 28(3), 165–179.
552 <https://doi.org/10.3109/03005369409086564>
- 553 Vance, M., Rosen, S., & Coleman, M. (2009). Assessing speech perception in young
554 children and relationships with language skills. *International Journal of*
555 *Audiology*, 48(10), 708–717. <https://doi.org/10.1080/14992020902930550>

556 Vickers, D. A., Moore, B. C. J., Majeed, A., Stephenson, N., Alferaih, H., Baer, T., &
557 Marriage, J. E. (2018). Closed-Set Speech Discrimination Tests for Assessing
558 Young Children. *Ear and Hearing*, 39(1), 32–41.
559 <https://doi.org/10.1097/AUD.0000000000000528>

560 Wightman, F. L., & Kistler, D. J. (2005). Informational masking of speech in children:
561 Effects of ipsilateral and contralateral distracters. *The Journal of the Acoustical*
562 *Society of America*, 118(5), 3164–3176. <https://doi.org/10.1121/1.2082567>

563 Ziegler, J. C., Pech-Georgel, C., George, F., Alario, F. X., & Lorenzi, C. (2005).
564 Deficits in speech perception predict language learning impairment.
565 *Proceedings of the National Academy of Sciences of the United States of*
566 *America*, 102(39), 14110–14115.
567 <http://www.pnas.org/content/102/39/14110.abstract>

568 Ziegler, J. C., Pech-Georgel, C., George, F., & Lorenzi, C. (2009). Speech-
569 perception-in-noise deficits in dyslexia. *Developmental Science*, 12(5), 732–745.
570 <https://doi.org/10.1111/j.1467-7687.2009.00817.x>

571 Zoubrinetzky, R., Collet, G., Serniclaes, W., Nguyen-Morel, M. A., & Valdois, S.
572 (2016). Relationships between categorical perception of phonemes, phoneme
573 awareness, and visual attention span in developmental dyslexia. *PLoS ONE*,
574 11(3), 1–26. <https://doi.org/10.1371/journal.pone.0151015>

575

576 **Supplementary materials**

577 S1: Full List of targets and Foils for the familiarisation and testing phase. AoA = Age
578 of Acquisition, SAM-PA = SAM-PA machine readable IPA transcription, Feature =
579 Phonological feature change, SNR = SNR adjustment for each word.

580

	Target X- SAMP A	AoA	SNR	Foil 1 (distracter) X- SAM PA	Feature	Foil 2 (distracter) X- SAMP A	Feature
Famili arisati on							
Bike	balk	4.79	2	walk	Manner	galk	Place
Bin	bln	4.68	3	mln	Manner	gln	Place
Bus	bVs	3.85	-4	wVs	Manner	dVs	Place
Dog	dQg	2.80	2	nQg	Manner	gQg	Place
Doll	dQl	3.68	0	rQl	Manner	bQl	Place
Duck	dVk	3.50	-4	zVk	Manner	gVk	Place
Laugh	lA:f	3.79	-2	zA:f	Manner	wA:f	Place
Leg	leg	3.00	-1	deg	Manner	jeg	Place
One	wVn	3.23	-3	mVn	Manner	lVn	Place
Rain	reln	3.60	0	neln	Manner	jeln	Place

Sea	si:	4.74	-7	zi:	Voicing	Ti:	Place
Sun	sVn	3.40	11	zVn	Voicing	TVn	Place
Watch	wQtS	4.33	-3	qQtS	Manner	rQtS	Place
Wave	welv	4.26	-1	belv	Manner	lelv	Place
Test items							
Bed	bed	2.89	-3	med	Manner	ped	Voicing
Book	bUk	3.68	0	wUk	Manner	pUk	Voicing
Boot	bu:t	3.89	5	wu:t	Manner	pu:t	Voicing
Chair	tSe@	3.43	0	se@	Manner	dZe@	Voicing
Boat	b@Ut	3.84	-1	w@Ut	Manner	p@Ut	Voicing
Bag	b{g	4.28	-3	m{g	Manner	p{g	Voicing
Dig	dig	4.19	-3	nlg	Manner	tlg	Voicing
Towel	taUl	3.22	-5	saUl	Manner	paUl	Place
Sing	sIN	3.47	-13	tIN	Manner	SIN	Place
Knife	nalf	4.15	0	dalf	Manner	malf	Place
Wash	wQS	4.00	-5	bQS	Manner	rQS	Place
Bath	bA:T	3.23	-4	wA:T	Manner	dA:T	Place
Leaf	li:f	4.60	2	ni:f	Manner	wi:f	Place
Road	r@Ud	4.55	-2	z@Ud	Manner	j@Ud	Place
Cough	kAf	4.32	18	pAf	Place	gAf	Voicing
Bite	balt	3.58	-5	dalt	Place	palt	Voicing
Comb	k@Um	5.50	9	p@U m	Place	g@Um	Voicing
Kite	kalt	4.58	5	palt	Place	galt	Voicing
Cow	kaU	3.94	0	taU	Place	gaU	Voicing
Cake	kelk	3.26	3	pelk	Place	gelk	Voicing
Fish	fIS	4.05	1	hIS	Place	vIS	Voicing
Fork	fO:k	3.63	4	sO:k	Place	vO:k	Voicing
Five	s@Up	4.51	4	Salv	Place	valv	Voicing
Fall	fO:l	4.71	0	sO:l	Place	vO:l	Voicing
Soap	s@Up	3.17	2	f@Up	Place	z@Up	Voicing
Foot	fUt	3.44	4	hUt	Place	vUt	Voicing
Suck	sVk	5.58	-8	hVk	Place	zVk	Voicing
Thumb	TVm	4.42	3	SVm	Place	DVm	Voicing

581