



**MACQUARIE**  
University  
SYDNEY · AUSTRALIA

## Macquarie University PURE Research Management System

---

**This is the accepted author manuscript version of an article published as:**

Xi, X., Li, J. N., Yuen, K. C., Chen, A. T., Li, S. Q., Hong, M. D., ... & Ching, T. Y. (2023). List Equivalency and Critical Differences of a Mandarin Bamford-Kowal-Bench Sentence in Babble Noise Test for Adults and Preschool Children With Normal Hearing. *Journal of Speech, Language, and Hearing Research*, 66(12), 5061-5070.

**Access to the published version:** [https://doi.org/10.1044/2023\\_JSLHR-23-00025](https://doi.org/10.1044/2023_JSLHR-23-00025)

**Copyright:** Version archived for private and non-commercial use with the permission of the author/s. For further rights please contact the author/s or copyright owner.

1 **List equivalency and critical differences of a Mandarin BKB sentence**  
2 **in babble noise test for adults and pre-school children with normal**  
3 **hearing**

4 Xin Xi<sup>a,b\*</sup>, Jia-Nan Li<sup>a,b\*</sup>, Kevin C. P. Yuen<sup>c</sup>, Ai-Ting Chen<sup>a,b</sup>, Si-Qi Li<sup>d</sup>,  
5 Meng-Di Hong<sup>a,b</sup>, Qian Wang<sup>a,b</sup>, Fei Ji<sup>a,b</sup>, Harvey Dillon<sup>e,f</sup> and Teresa Y. C.  
6 Ching<sup>g, h, i</sup>.

7 \*These authors have equal contribution to this paper and the test development.

8

9 *<sup>a</sup>Department of Otolaryngology - Head & Neck Surgery, the Sixth Medical Center of PLA*  
10 *General Hospital, Beijing, China;*

11 *<sup>b</sup>National Clinical Research Center for Otolaryngologic Diseases, Beijing, China;*

12 *<sup>c</sup>Department of Special Education and Counselling, The Education University of Hong*  
13 *Kong, HKSAR, China;*

14 *<sup>d</sup>School of Communication Science, Beijing Language and Culture University, Beijing,*  
15 *China;*

16 *<sup>e</sup>Department of Linguistics, Macquarie University, Sydney, Australia;*

17 *<sup>f</sup>Manchester Centre for Audiology and Deafness, School of Health Sciences, University*  
18 *of Manchester, United Kingdom;*

19 *<sup>g</sup>Macquarie School of Education, Macquarie University, Sydney, Australia;*

20 *<sup>h</sup>NextSense Institute, Sydney, Australia;*

21 *<sup>i</sup>School of Health and Rehabilitation Sciences, University of Queensland, Brisbane,*  
22 *Australia.*

23

24 Correspondence:

25 Teresa Y. C. Ching, Macquarie School of Education, Macquarie University, Sydney  
26 2109, Australia; [teresayching@gmail.com](mailto:teresayching@gmail.com)

27 Harvey Dillon, Department of Linguistics, Macquarie University, Sydney 2109,  
28 Australia; [harvey.dillon@mq.edu.au](mailto:harvey.dillon@mq.edu.au)

29

30 **Abstract**

31 **Purpose:** To determine the speech recognition equivalence of Mandarin BKB sentence  
32 lists with adults and children with normal hearing.

33 **Method:** A total of 32 lists each of 9 sentences were compiled from a corpus of BKB-  
34 like sentences with paired babble in Mandarin. Inter-list equivalence, critical  
35 differences, and sensitivity of performance to signal-to-noise ratio (SNR) were  
36 examined. Experiment 1 included 64 native Mandarin-speaking adults with normal  
37 hearing. Experiment 2 included 54 native Mandarin-speaking children with normal  
38 hearing aged 4–6 years old.

39 **Results:** Among the 32 sentence lists, 28 lists were confirmed to be equivalent in  
40 adults, with a mean SNR required for 50% correct (SNR<sub>50</sub>) of  $-5.9 \pm 0.1$  dB, a mean  
41 slope of  $22.3 \pm 1.5\%/dB$ , and the grand 95% critical difference (CD) was subsequently  
42 calculated as 27.2% for score. From the 28 equivalent lists, 27 lists were selected and  
43 observed to be equivalent in children, with a mean SNR<sub>50</sub> threshold of  $-2.0 \pm 0.2$  dB, a  
44 mean slope of  $15.8 \pm 1.1\%/dB$ , and a grand 95% CD of 24.6% for score.

45 **Conclusion:** The Mandarin BKB sentences in babble noise test offers an opportunity  
46 for clinicians and researchers to assess speech understanding in adults and pre-school  
47 children in an efficient manner. For comparisons of performance in different test  
48 conditions, 28 equivalent lists are available for adults, and 27 equivalent lists for pre-  
49 school children. The 95% critical difference values can be used for total percentage  
50 correct or SNR for 50% performance. Future work will examine the clinical utility for  
51 school-aged children and children who are deaf and hard of hearing.

52

53 **Keywords:** speech perception; Mandarin sentence test; babble noise

54 **Introduction**

55           Speech audiometry is a fundamental tool in the assessment of hearing disorders  
56 and the evaluation of rehabilitative outcomes. It can not only aid in determining the  
57 degree and type of hearing loss in conjunction with pure-tone audiometry, but also help  
58 make clinical decisions about appropriate gain and maximum output of hearing aids and  
59 other auditory prosthetics for hearing-impaired patients, as well as help assess how well  
60 they hear in noise.

61           The primary concern of most hearing-impaired people is the difficulties in  
62 understanding speech in noisy environments (Moore, 2003), which cannot be accurately  
63 predicted from hearing thresholds (Killion et al., 2004). Therefore, speech-in-noise tests  
64 have been recommended as routine practice for clinical evaluation of patients using  
65 hearing aids (HAs) and cochlear implants (CIs) (Wilson, 2004). The target speech can  
66 range in complexity from nonsense syllables, digits, to monosyllabic words,  
67 multisyllabic words, and sentences. Compared with word recognition tests, sentence  
68 materials contain redundant information and dynamic amplitude variation that are  
69 typical of connected speech. Given its better resemblance to speech situations that users  
70 of HAs or CIs would encounter in realistic listening environments, sentence tests may  
71 be a more valid measurement of real-world speech intelligibility.

72           Since HAs and CIs have become increasingly common for the intervention of  
73 hearing loss in China, there has been an urgent need for standardized speech recognition  
74 test materials. As the official language of China, Mandarin Chinese is spoken by over  
75 80% of its population (China State Language Commission, 2021). The popularity of  
76 Mandarin enables the measurement of speech perception among people from different  
77 dialect regions. Currently, there are a limited number of sentence recognition tests  
78 available in Mandarin Chinese. The first standardized sentence test is the Mandarin  
79 Hearing in Noise Test (MHINT) developed by Wong and her colleagues (2007)

80 following the same rationale as the English HINT (Nilsson, Soli & Sullivan, 1994). The  
81 MHINT comprises 24 lists of 10 sentences each, with normative data available for  
82 measuring the reception threshold for sentences in quiet as well as in noise. Fu, Zhu and  
83 Wang (2011) focused on the differences in tone perception between listeners with  
84 normal hearing and CI users and developed a set of phonetically balanced sentence lists  
85 named the Mandarin speech perception (MSP) test. Hu et al. (2018) developed a  
86 sentence test using semantically unpredictable sentences with fixed grammatical  
87 structure to assess speech perception in noise, namely the Mandarin Chinese matrix  
88 (CMN matrix) sentence test.

89         When it comes to pediatric sentence tests, a closed-set sentence identification  
90 test (Mandarin Paediatric Speech Intelligibility test, MPSI; Zheng et al., 2009) has been  
91 developed to evaluate children’s speech recognition first in quiet and then in the  
92 presence of a competing sentence at signal-to-noise ratios (SNRs) of +10, +5, 0, –5 and  
93 –10 dB. Since there are children scoring at the ceiling of the test, Chen and Wong  
94 (2020) created a sentence test with a higher level of difficulty. They selected sentences  
95 suitable for evaluating children aged 6 years and above (primary school age in Mainland  
96 of China) from the MHINT, and developed an open-set pediatric sentence test – the  
97 Mandarin version of the Hearing in Noise Test for Children (MHINT-C). The MHINT-  
98 C comprises 12 lists each of 10 sentences, with each sentence consisting of 10  
99 characters. The test was verified for use with Mandarin-speaking children aged 6–18  
100 years with age-specific correction factors established. The ability to understand speech  
101 in the presence of competing sounds, particularly speech, is especially important for  
102 preschool aged children, who spend much of their waking hours in acoustic  
103 environments where noise (speech or otherwise) levels can be high while they learn and  
104 acquire language (Newman et al., 2015; Leibold & Buss, 2019). Further work is needed

105 to examine how pre-school children perform in environments such as classrooms, and to  
106 devise appropriate tools for evaluating young children's performance using hearing aids  
107 and cochlear implants (Litovsky et al., 2006). Given the trend of an increasing number  
108 of young candidates for cochlear implantation (Li et al., 2017), it is necessary to devise  
109 open-set sentence tests for the evaluation of speech perception in noise abilities in  
110 children younger than 6 years of age (Ching et al, 2018) to be used in addition to  
111 closed-set sentence tests.

112 In English-speaking countries, the Bamford-Kowal-Bench (BKB) sentences  
113 (Bench & Bamford, 1979) are one of the most commonly used speech materials. BKB  
114 sentences are short, highly redundant and rich in semantic and syntactic context. So,  
115 they are appropriate for assessing the speech recognition of children as well as adults.  
116 The BKB sentence test is presented in an open-set format and is commonly used for  
117 candidacy selection for CIs, according to FDA-approved criteria. Several speech-in-  
118 noise test materials have also been developed based on BKB sentences, such as the  
119 Hearing in Noise Test (HINT; Nilsson, Soli & Sullivan, 1994) and the BKB-Sentences  
120 in Noise test (BKB-SIN; Etymotic Research, 2005). The BKB-SIN test has been  
121 included in the latest version of the Minimum Speech Test Battery for adult CI users  
122 (MSTB, 2011).

123 Given the effectiveness and practicality of BKB sentences, they can be a proper  
124 reference for standardized Mandarin sentence tests. In line with the principles of the  
125 BKB sentence test, Xi and his colleagues (2012) constructed a corpus of Mandarin  
126 BKB-like sentences in babble noise with *Homogeneity Optimized via Psychometric*  
127 *Evaluation*, the HOPE corpus. In the study, 647 6- to 8-character sentences within the  
128 syntactic and semantic abilities of children were selected from a speech corpus of pre-  
129 school children (Liu et al. 2008). A total of 309 sentences were retained based on the

130 speech-in-noise recognition performance of 96 Mandarin-speaking adults. Following  
131 the International Collegium of Rehabilitative Audiology (ICRA) guidelines for the  
132 development of multilingual speech tests (Akeroyd et al., 2015), the SNR of each  
133 keyword was adjusted by the difference between that keyword's speech recognition  
134 threshold ( $SNR_{50}$ ) and the average  $SNR_{50}$  across all 309 sentences, to equalize the  
135 difficulty of test items. The homogeneity of adjusted keywords was verified in another  
136 group of 64 Mandarin-speaking subjects. Although the design of the test material is  
137 appealing for hearing evaluation in young children and research, the psychometric data  
138 available on the sentences, to date, were based on adult performance. As the sentences  
139 were not compiled into lists, information on list equivalency and critical difference were  
140 not available.

141         The purpose of the current study was to compile a set of Mandarin sentence in  
142 babble noise test lists from the homogeneous sentence-babble pairs in the HOPE corpus  
143 and determine the list equivalency and critical differences for adults and pre-school  
144 children with normal hearing.

## 145 **Experiment 1: List construction and equivalency in adults with normal** 146 **hearing**

### 147 *Materials and methods*

#### 148 *Materials*

149         A total of 32 sentence lists (Lists A to Z and Lists 1 to 6), each comprising nine  
150 scoring sentences and one preliminary example sentence with paired frozen babble  
151 noise, were constructed from the HOPE corpus. The target sentences were read by a 20-  
152 year-old male native Mandarin speaker, while the babble noise was created from



153 connected discourse read by two males and two females aged from 28 to 38 (Details of  
154 the generation of target speech and babble noise can be seen in Xi et al., 2012).

155         The detailed procedures of list compilation were described as follows. First, all  
156 309 sentences in the corpus were categorized in terms of sentence type (statements,  
157 interrogatives, exclamatory, and imperatives) and ranked according to syntactic  
158 difficulty corresponding to the developmental stages of children. According to Zhu and  
159 Miao's (1990) theory of child language development, syntactic difficulty of sentences  
160 can be categorized into four types – simple sentence without any modifiers, complex  
161 predicate sentence, complex subject/object sentence and compound sentence. Syntactic  
162 development progressed from acquisition of simple sentences without any modifiers to  
163 compound sentences. Next, 288 sentences were selected as scoring sentences to form 32  
164 lists based on three rules: (1) sentences belonging to different types and syntactic  
165 difficulty should be balanced across lists as much as possible; (2) sentences about the  
166 same topic should be distributed as balance as possible across lists; for example,  
167 sentences talking about having food would not appear too many times in a single list;  
168 (3) the number of keywords for scoring in each list was 50. The remaining sentences  
169 from the corpus were then used as the preliminary example sentence in each list, and  
170 also in four practice lists.

171         The keyword-based scoring rule and selection of keywords remained the same as  
172 the study on corpus construction (Xi et al., 2012) since item homogeneity was verified  
173 in terms of keywords. All keywords for scoring were selected from content words and  
174 function words that received primary stress in each sentence. Each keyword contained  
175 one to three characters. The unit of measurement was defined as word instead of  
176 character because single characters in Chinese can be combined to form a word  
177 representing a semantic unit. With additional semantic cues, listeners tend to recognize

178 the word as a whole rather than separated single characters. The number of keywords  
179 for scoring in each sentence ranged from three to seven, with a total of 50 keywords in  
180 each list.

181 All 32 test lists and the four practice lists were burned onto a compact disc (CD).  
182 Each sentence, in one channel, was preceded by a 1.6 kHz beep 500 ms before the start  
183 of the sentence and was time-locked to a paired segment of babble noise in another  
184 channel. The babble noise commenced 2 seconds prior to each sentence and ceased 500  
185 ms after the sentence. The inter-stimuli interval was four seconds. The CD also  
186 contained a 60-second calibration noise with its spectrum matched with International  
187 Long-Term Average Speech Spectrum (ILTASS; Byrne et al, 1994), which  
188 corresponded with the stimulus sentences and babble segments that had also been  
189 filtered to match the ILTASS (Xi et al., 2012). The calibration noise's level was set at  
190 the same RMS (Root of Mean Square) of equivalent continuous level as the sentences  
191 and babble segments.

## 192 *Participants*

193 A total of 64 native Mandarin-speaking adults (42 females and 22 males) with no  
194 known history of hearing or speech disorders were recruited to participate. All  
195 participants were born and grew up in Beijing or nearby provinces so that the effects of  
196 dialect were minimized. They were aged 18–30 years (mean = 23.8 years, SD = 2.8  
197 years), with air-conduction pure-tone thresholds  $\leq 15$  dB HL at audiometric frequencies  
198 from 250 to 8000 Hz as well as normal type A tympanograms. The recruitment of  
199 human subjects in the present study was reviewed and approved by the Institutional  
200 Review Board of the Chinese PLA General Hospital, Beijing, China.

201 *Experimental design and procedure*

202 The experiment was conducted in a sound-proof booth in the Chinese PLA General  
203 Hospital, Beijing. The sentences were presented via a Marantz 5400 CD player to  
204 channel 1 of a GSI 16 audiometer at a fixed level of 60 dB SPL. The babble noise was  
205 presented to Channel 2 of the audiometer, with its intensity manually adjusted by the  
206 experimenter. As the SNR<sub>50</sub> for adult listeners with normal hearing was found to be  
207 approximately -6 dB SNR (SD = 1.1 dB) in the HOPE corpus study (Xi et al., 2012),  
208 sentences were tested at -9, -7, -5 and -3 dB SNRs in this experiment. The sentences  
209 (Channel 1) and babble noise (Channel 2) were mixed by the audiometer and connected  
210 to a GSI loudspeaker 1 m in front of the participants. The sound field was calibrated  
211 according to the Sound Field Measurement Tutorial (American Speech-Language-  
212 Hearing Association, 1991). A Brüel & Kjaer 2250 sound level meter was used to  
213 measure the loudspeaker's output of the calibration noise with its microphone located at  
214 the reference point corresponding to the position of the participant's head. Level  
215 adjustments for sentences and noise were based on the measured level of calibration  
216 signal.

217 The participants were instructed to repeat as many words in the sentences as  
218 possible after each sentence was presented. Performance was scored by the  
219 experimenter during the test in terms of number of keywords correctly repeated. Four  
220 practice lists were presented at all testing SNR levels (i.e., -9, -7, -5 and -3 dB SNRs) in  
221 a random order to acquaint participants with the listening environment and the listening  
222 task before testing. The test sessions last about 2 hours for adult participants with rest  
223 included. The experimenter monitored participants' fatigue and attention, and offered  
224 breaks when necessary.

225 The 32 sentence lists were presented at the four SNRs with both test order and

226 the presentation SNR counterbalanced across subjects (see Supplemental Material 1). In  
227 this way, each sentence list was tested by 16 participants at each SNR level. Learning  
228 and/or fatigue effects were minimized as the testing order of all lists were  
229 counterbalanced across participants.

## 230 ***Results***

### 231 *Analysis of list equivalence*

232 The performance scores of all lists were calculated as the proportion of correct  
233 recognition of the 50 keywords for a particular list at each SNR averaged across 16  
234 participants (see Supplemental Material 2 for descriptive data). A two-way repeated  
235 measures analysis of variance (ANOVA) was performed on the scores with SNR level  
236 as a repeated-measure factor to investigate the equivalence of the 32 sentence lists.  
237 Significant main effects of list number ( $F_{31,465} = 3.953, p < 0.001$ ) and SNR ( $F_{3,45} =$   
238  $5201.672, p < 0.001$ ) were observed. Post hoc Tukey Honestly Significant Difference  
239 (HSD) test revealed that the recognition scores in List K, List Q, List U and List Y  
240 differed significantly from other lists. The four non-equivalent lists were subsequently  
241 removed from further analysis.

### 242 *Performance-SNR function*

243 The  $SNR_{50}$  and slope of performance-SNR intelligibility functions for each remaining  
244 sentence list were obtained by fitting a logistic regression function using the following  
245 equation (Nissen et al., 2005):

$$246 \quad p = \frac{\exp(a - b * SNR)}{1 + \exp(a - b * SNR)} \quad (1)$$

247 in which the dependent variable  $p$  (percentage correct) is the performance score,  $a$  is the  
248 regression intercept, and  $b$  is the regression slope. The slope (percentage per dB) of the

249 steepest portion of the curve for each list was calculated as  $-b/4$ . The speech recognition  
250 threshold ( $SNR_{50}$ ) was calculated as the ratio of  $a/b$ .

251 For the 28 sentence lists, the mean and standard deviation of  $SNR_{50}$  thresholds  
252 were  $-5.9 \pm 0.1$  dB, and the slope of performance-SNR functions,  $22.3 \pm 1.5\%/dB$ .  
253 Psychometric function curves of all lists are shown in Figure 1.

254 [Figure 1]

### 255 *Variability of scores*

256 All recognition scores in percentage were normalized for comparison using the  
257 rationalized arcsine unit (rau) transform (Studebaker, 1985; Sherbecoe & Studebaker,  
258 2004). This transformation can achieve a linear and additive scale, and come up with  
259 rationalized arcsine values nearly equal to percentages over a broad central range of  
260 percentages. The scores of the 28 lists at four SNR conditions were transformed via the  
261 following equations:

$$262 \quad \theta = \arcsin \sqrt{\frac{X}{N+1}} + \arcsin \sqrt{\frac{X+1}{N+1}} \quad (2)$$

$$263 \quad R = \left(\frac{146}{\pi}\right) \theta - 23 \quad (3)$$

264 in which  $X$  represents the number of correct items, and  $N$  is the total number of tested  
265 items.

266 The standard deviations of normalized scores of each participant at each SNR  
267 were calculated. The grand SD was computed by averaging the squared standard  
268 deviations and taking the square root, resulting in a mean value of 9.8%. Given the  
269 slope of the P-SNR function (22.3%/dB), and the symmetrical distribution of scores  
270 around the 50% point, the corresponding inter-participant SD of  $SNR_{50}$  was 0.4 dB. By  
271 multiplying the grand SD by a coefficient of 2.77 (derived from  $\sqrt{2} \times 1.96$  where 1.96 is

272 the appropriate  $z$  value for 95% confidence interval, and  $\sqrt{2}$  indicates that the critical  
273 difference includes variation in two measurements), the 95% critical difference (CD)  
274 was subsequently calculated as 27.2% for score and 1.2 dB SNR for speech recognition  
275 threshold ( $\text{SNR}_{50}$ ).

## 276 **Experiment 2: List equivalency in pre-school children with normal hearing**

### 277 *Materials and methods*

#### 278 *Materials*

279 The 28 sentence lists observed to be equivalent in adults with normal hearing in  
280 Experiment 1 were administered to children with normal hearing. Since three SNR  
281 conditions were determined to be tested in pediatric participants, the total number of  
282 lists need to be adjusted to a multiple of three so that all lists can be counterbalanced  
283 across SNR conditions. Therefore, one list (List 6, with the greatest variance from other  
284 lists though insignificant) was removed, resulting in a total of 27 sentence lists tested in  
285 children with normal hearing. Each list consisted of 9 scoring sentences with 50  
286 keywords, and a preliminary example sentence.

#### 287 *Participants*

288 A total of 54 native Mandarin-speaking children (27 male, 27 female) aged 4–6  
289 years were recruited from a kindergarten in Beijing to participate in Experiment 2. To  
290 balance the age of the pediatric participants, 18 children were recruited for each of three  
291 age group: 4 to 4.5 years old, 4.5 to 5 years old, and 5 to 6 years old. All children had  
292 no known history of otologic diseases, with pure tone thresholds for each audiometric  
293 frequency from 250 Hz to 8 kHz  $\leq 25$  dB HL and normal results of otoscopic  
294 examination and tympanometry. The recruitment of human subjects in the present study

295 was reviewed and approved by the Institutional Review Board of the Chinese PLA  
296 General Hospital, Beijing, China.

297 *Experimental design and procedure*

298 The pediatric participants attended two sessions because of the children's short attention  
299 spans and low cooperative motivation. In the first session, otoscopic examination, pure  
300 tone audiometry and tympanometry were administered to screen the pediatric  
301 participants, then the suitable children were familiarized with the speech test procedure.  
302 In the second session, the formal sentence recognition test was administered after at  
303 least two practice lists were completed. The children were instructed to listen carefully  
304 to each sentence after the beep and repeat as many words as possible. Performance was  
305 scored by the experimenter during the test in terms of number of keywords correctly  
306 repeated. The test session lasted an average of 1.5 hours with break periods included.  
307 The experimenter monitored children's fatigue and attention, and offered verbal  
308 encouragement, food reward, or breaks when necessary to encourage the children to  
309 complete all required testing.

310         The test was administered in a quiet room in a kindergarten. The ambient noise  
311 was  $\leq 40$  dB (A). The instruments (CD player, audiometer, and loudspeaker) used for  
312 the presentation of stimulus sentences and babble noise were the same as in Experiment  
313 1. The sentences were presented at a fixed level of 65 dB SPL, slightly higher than the  
314 level in Experiment 1 given the higher level of ambient noise. The intensity of the  
315 babble noise was adjusted by the experimenter to achieve +1, -2, -5 dB SNRs. The SNR  
316 levels were pre-determined based on the results of a pilot test with six children. They  
317 were first tested at the highest SNR level used in the experiment with adults, i.e., -3 dB  
318 SNR, and scored around 35%. We then increased the SNR to -2 dB SNR, and their  
319 performance scores reached about 50%, suggesting that -2 dB SNR would produce

320 approximately 50% intelligibility for pediatric listeners with normal hearing. The  
321 loudspeaker was placed 1 m in front of the participants and approximately 30 cm from  
322 one corner of the room with its center at the same height as the ear of the child in a  
323 sitting position. The sound field was calibrated in the same way as in Experiment 1.

324         Considering children's limited attention span, only 15 of the 27 lists were tested  
325 for each child. The sentence lists were assigned in a balanced way to three age groups in  
326 the following way: first, children were divided into three age groups (aged 4-4;6, aged  
327 4;6-5, and aged 5-6) with 18 participants in each group; then each list were distributed  
328 to be tested 10 times in each group ( $18 \text{ children} \times 15 \text{ lists} = 27 \text{ lists} \times 10 \text{ times}$ ). The 15  
329 sentence lists were presented at different SNRs, with the starting list varying across  
330 subjects (see Supplemental Material 3). In this way, all lists were optimally  
331 counterbalanced across SNRs resulting in each list tested by 10 different children in  
332 each of three SNR conditions. Learning and/or fatigue effects were minimized as the  
333 presentation order of the lists was also counterbalanced among the participants, with  
334 each list having equal probability of being presented first.

## 335 ***Results***

### 336 *Analysis of list equivalence*

337 The performance scores were calculated as the proportion of correct recognition of the  
338 50 keywords for each list at each SNR averaged across 10 participants (see  
339 Supplemental Material 4 for descriptive data). A two-way repeated measures ANOVA  
340 was performed using SNR level as a repeated-measure factor to investigate the  
341 equivalence of the 27 lists in pediatric participants. Significant main effects of list  
342 number ( $F_{26,234} = 1.86, p = 0.009$ ) and SNR ( $F_{2,18} = 6197.14, p < 0.001$ ) were observed  
343 for recognition scores. However, the follow-up Tukey's HSD post hoc test revealed no



344 significant differences among the 27 lists.

### 345 *Performance-SNR function*

346 The Performance-SNR psychometric functions for the 27 lists were obtained based on  
347 the same mathematical procedures as in Experiment 1. Logistic curves were fitted for  
348 the recognition scores of each list.

349 For the 27 lists, the mean and standard deviation of the slopes at the SNR<sub>50</sub>  
350 across the lists were  $15.8 \pm 1.1\%/dB$ , and the SNR<sub>50</sub>,  $-2.0 \pm 0.2$  dB. Psychometric  
351 function curves of all lists are shown in Figure 2.

352 [Figure 2]

### 353 *Variability of scores*

354 All recognition scores of the 27 lists at three SNR conditions ( $n = 10$ ) were normalized  
355 by a rationalized arcsine transformation as in Experiment 1. The grand SD was  
356 computed as 8.9% for score and 0.6 dB for SNR<sub>50</sub>. The critical difference was  
357 calculated in the same way as in Experiment 1, resulting in a 95% CD of 24.6% for  
358 score and 1.6 dB for SNR<sub>50</sub>.

### 359 *Age effect*

360 A simple linear regression was conducted to assess the effect of age on recognition  
361 scores. For SNR = -5 dB, the regression equation was:  $\text{score} = 0.004 + 0.025 \times \text{age}$ ,  
362 which was found to be statistically significant ( $R^2 = 0.046$ ,  $F_{1,268} = 12.944$ ,  $p = .000$ ).  
363 For SNR = -2 dB, the regression equation was:  $\text{score} = 0.112 + 0.08 \times \text{age}$ , which was  
364 found significant ( $R^2 = 0.171$ ,  $F_{1,268} = 55.349$ ,  $p = .000$ ). For SNR = +1 dB, the  
365 regression equation was:  $\text{score} = 0.557 + 0.063 \times \text{age}$ , which was found significant ( $R^2$

366 = 0.192,  $F_{1,268} = 63.667$ ,  $p = .000$ ). Generally, children's age significantly predict their  
367 recognition scores in all three SNR conditions.

## 368 **Discussion**

369 This paper provides evidence on list equivalency and critical differences for the  
370 Mandarin BKB sentence in noise test for adults and young children with normal  
371 hearing. Sentences with paired babble from the BKB-like sentence corpus in Mandarin  
372 Chinese (HOPE corpus, Xi et al., 2012) were compiled into 32 lists with balanced  
373 distributions of syntactic difficulty and sentence type. For the 28 lists that were found to  
374 be equivalent in difficulty for adults, a critical difference of 1.2 dB for  $SNR_{50}$  was  
375 obtained at the speech presentation level of 60 dB SPL. Of the 28 lists, 27 were  
376 equivalent in difficulty for children, and a critical difference of 1.6 dB for  $SNR_{50}$  was  
377 obtained for the speech presentation level at 65 dB SPL. Using percent correct scores  
378 for the equivalent lists, the critical difference was 27.2% for adults and 24.6% for  
379 children. These values suggest that any greater difference can be interpreted as  
380 significant at the 95% confidence level.

381 By establishing list equivalency and critical differences, the present study  
382 provides a valuable resource for future applications. In clinical practice, speech  
383 perception tests usually include several equivalent lists for repeated measurements to  
384 control for potential learning effects. When different lists are administered at different  
385 intensity levels, one can estimate the speech intelligibility curve, provided that the lists  
386 are interchangeable. That is, any listener receiving a particular percentage correct score  
387 on one list under certain listening conditions should receive the same score (within  
388 certain limits) on any other list if all other variables remain constant.

389 In clinical trials, lists that are not equivalent in difficulty can be used to compare  
390 scores between two conditions, e.g., for within-subject assessment of the effects of

391 different forms of hearing prostheses (e.g., HAs or CIs) or training. The difference in  
392 scores between two presentations must be large enough to be attributable to the  
393 experimentally varied condition rather than to chance alone. The smallest difference  
394 that must be exceeded for a difference between two scores is termed the critical  
395 difference. Considering the experimental design of the current study, this variability was  
396 quantified with inter- and intra-subject differences combined by calculating the standard  
397 deviation of the resulting distributions of scores in a group of subjects (Dillon, 1982). In  
398 this study, list scores for each listener were normalized before the calculation of critical  
399 difference using the rationalized arcsine unit (rau) transform to achieve a linear and  
400 additive scale. This transformation was proposed by Studerbaker (1985) for the data  
401 analysis of speech recognition scores. The purposes of the rau transform are firstly so  
402 that the random error component more closely approximates a Gaussian distribution  
403 even when the score is close to 0% or 100%, secondly to minimize the correlation of  
404 variability and performance, and thirdly so that the transformed scores approximate the  
405 original proportion correct over most of the performance range. The mean values of  
406 9.8% in adults and 8.9% in children represent the variability attributable to a  
407 combination of the differences in true list difficulty and random measurement error  
408 within a subject. The critical difference for two scores based on single lists to be  
409 significantly different, with a 95% confidence interval, is thus 27.2% for adults and  
410 24.6% for children. The critical difference is useful when evaluating the efficacy of a  
411 treatment or rehabilitation. These values approximate the estimate for the Australian  
412 version of BKB sentence lists (24.6%), which also contained 50 keywords in each list  
413 (Keidser et al., 2002). In both cases, the standard deviation of repeated scores (e.g.,  
414 9.8% for adults in this study) is slightly larger than the value of 7.1% that would be  
415 expected from the binomial distribution for a true score 50% correct in a sample of 50

416 independent items (Hagerman, 1976). The larger value observed is undoubtedly because  
417 probabilities of correctly recognizing the keywords in a semantically rich sentence are  
418 not statistically independent of each other, so effectively there are fewer than 50  
419 *independent* scored items per list.

420 Another significant issue that needs to be addressed in the design of standardized  
421 speech perception tests is sensitivity (Dillon & Ching, 1995). For a test to be sensitive to  
422 some specific change in listening conditions (e.g., SNR), the resulting change in scores  
423 must be larger than the variation expected on the basis of chance alone. Test materials  
424 with a steep slope of the P-SNR function are more sensitive than materials with  
425 shallower slopes because any small change in SNR is more likely to lead to a change  
426 larger than occur due to chance alone. Therefore, the steep slope of the P-SNR function  
427 is generally one of the focuses in the development of speech-in-noise tests. When  
428 testing the sentence lists in adults and children with normal hearing in the current study,  
429 the average slopes of P-SNR functions were  $22.3 \pm 1.5\%/dB$  (28 lists) and  $15.8 \pm$   
430  $1.1\%/dB$  (27 lists), respectively. Compared with the slope of  $10\%/dB$  in the HINT  
431 (Nilsson et al., 1994) and  $9\%/dB$  in the MHINT (Wong et al., 2007), the current set of  
432 sentence-recognition-in-babble-noise lists has a steeper slope of the P-SNR function.  
433 This difference may be partly attributed to the use of steady-state competing noise in the  
434 HINT and MHINT, and partly to the relative low difficulty of the test material as the  
435 tests were designed for school-aged children. Test sensitivity and the slope of the  
436 psychometric function are maximized by making all items in the test equally difficult  
437 (Dillon, 1983). For our Mandarin BKB sentence test, item homogeneity was achieved  
438 through level adjustment of each keyword for scoring so that all words had the same  
439  $SNR_{50}$  (Xi et al., 2012) following the ICRA recommendations for the construction of  
440 multilingual tests (Akeroyd et al., 2015). Such optimization based on word-specific

441 SNR<sub>50</sub> has also been applied in the development of CMN matrix, resulting in a  
442 relatively steep mean slope of 13.1%/dB (Hu et al., 2018). The results are consistent  
443 with previous studies showing that the more homogeneous performance is on the  
444 individual test items with respect to both location and slope, the steeper the slope of the  
445 mean psychometric function (Wilson & Carter, 2001). This keyword-based  
446 homogenization method contrasts with that used in the development of the MHINT test,  
447 where sentence, rather than individual word, difficulty was equalized, suggesting that  
448 the word-specific optimization of homogeneity can result in better sensitivity.  
449 Accordingly, the current sentence test has higher sensitivity to differences between test  
450 conditions and can thus evaluate different algorithms and settings of hearing prostheses  
451 with respect to speech intelligibility in noisy environments with better accuracy and  
452 efficiency (Wilson & Carter, 2001).

453         The strength of this study lies in establishing the list equivalency of 28 lists of  
454 BKB like Mandarin sentences drawn from a corpus of sentences each with a paired  
455 babble noise in which the keywords for scoring are homogeneity optimised. The 28 lists  
456 can be used with adults, and 27 of them also with pre-school children. Although list  
457 equivalency and critical differences have been established, this study does not provide  
458 evidence for clinical applications with adults and children who are deaf or hard of  
459 hearing. Future research needs to be carried out to investigate the clinical applicability  
460 of the lists in evaluating performance with hearing prostheses, including hearing aids  
461 and cochlear implants.

462         Unlike most existing speech perception tests in Mandarin Chinese that use  
463 steady-state speech-shaped noise as competing sounds, the present Mandarin BKB lists  
464 used 4-talker babble as competing sounds so that the test has greater face validity in  
465 measuring an individual's ability to understand speech in everyday acoustic

466 environments for speech communication. The use of a sentence paired with a frozen  
467 babble allowed effective optimisation of homogeneity across keywords used for scoring  
468 (Xi et al, 2012). The 4-talker babble provides energetic as well as informational  
469 masking, while the sentences are natural utterances based on a corpus of utterances of  
470 children aged between 4 and 5 years (Liu et al., 2008). Accommodating the cognitive,  
471 motor, and attentional capabilities of young children, the sentence length was controlled  
472 to 6–8 characters, shorter than sentences in MHINT (10 characters; Wong et al., 2007)  
473 and CMN matrix (11 characters; Hu et al., 2018). A single list takes approximately 1.5  
474 min to administer, which is appropriate for clinical use among children 4 years of age  
475 and older.

476         Although we have attempted to counterbalance lists and SNRs across  
477 participants in the study using a Latin square design (Supplementary Material 3), the  
478 order of SNRs was constant across participants during test administration. To reduce the  
479 potential effect of order on performance, each participant was tested at each SNR for  
480 five times such that testing at a certain SNR occurred sometimes at the beginning and  
481 other times towards the end of a test session. The potential effect of the order of SNRs  
482 will need to be quantified in future work.

483         The effect of age on performance observed in this study has been reported in  
484 other studies (e.g., the MHINT-C (Chen & Wong, 2020), the Cantonese Hearing in  
485 Noise Test for Children (CHINT-C; Wong, Chen & Leung, 2019), the Canadian French  
486 HINT for children (Vaillancourt et al., 2008). Future research will attempt to replicate  
487 the present findings on list equivalency and examine test-retest reliability in children  
488 across a wide age range.

489 **Conclusions**

490 The Mandarin BKB sentences in noise test offers an opportunity for clinicians and  
491 researchers to assess speech understanding in adults and pre-school children in an  
492 efficient manner. For clinical comparisons of performance with different presentation  
493 conditions or settings, 28 equivalent lists are available for adults, and 27 lists for  
494 preschool children. These lists can be scored using keyword scoring. Furthermore, the  
495 95% critical difference values can be used for total percentage correct or SNR for 50%  
496 performance to assess the effectiveness of listening in real-life environments. Future  
497 work on this research tool should examine clinical applicability for school-aged children  
498 and children who are deaf and hard of hearing.

499

500 **Author Note**

501 Portions of this article were previously described in Chinese in the following article and  
502 permission for reuse has been obtained from the publisher:

503 Xi, X., Chen, A. T., Li, J. N., Ji, F., Hong, M. D., Yang, S. M., Han, D.Y. (2009).

504 Caozayu zaosheng xia Putonghua ertong yuju ceting biao de biao zhunhua [Standardized  
505 Mandarin Sentence Perception in Babble Noise Test Materials for Children]. *Journal of*  
506 *Audiology and Speech Pathology*, 17(4), 318-322.

507

508 **Acknowledgments**

509 The authors acknowledge Mr. Yang Zhao's efforts in adults' data collection and the  
510 support of the NIHR Manchester Biomedical Research Centre.

511

512 **Funding**

513 This research is sponsored by National Key R&D Program of China (Grant No.  
514 2020YFC2004005, 2007BAI18B12, 2008BAI50B01), and supported by Science  
515 Foundation of Beijing Language and Culture University (the Fundamental Research  
516 Funds for the Central Universities, 21PT01).

517

518 **ORCID ID**

519 Xin Xi <https://orcid.org/0000-0001-7806-8396>

520 Siqi Li <https://orcid.org/0000-0001-8660-6414>

521

522 **Disclosure Statement**

523 The authors declared no potential conflicts of interest with respect to the research,  
524 authorship, and/or publication of this article.

525

526 **Data Availability Statement**

527 The datasets generated and analyzed during the current study, as well as the recorded  
528 speech materials and instructions for test administration, are available from the first  
529 author, Dr. Xin Xi via email [xixin\\_plagh@yeah.net](mailto:xixin_plagh@yeah.net) upon reasonable request.



530 **References:**

- 531 Akeroyd, M. A., Arlinger, S., Bentler, R. A., Boothroyd, A., Dillier, N., Dreschler, W.  
532 A., Gagné, J. P., Lutman, M., Wouters, J., Wong, L., Kollmeier, B., &  
533 International Collegium of Rehabilitative Audiology Working Group on  
534 Multilingual Speech Tests. (2015). International Collegium of Rehabilitative  
535 Audiology (ICRA) recommendations for the construction of multilingual speech  
536 tests. ICRA Working Group on Multilingual Speech Tests. *International*  
537 *Journal of Audiology*, 54 Suppl 2, 17–22.
- 538 American Speech-Language-Hearing Association. (1991). Sound field measurement  
539 tutorial. *Asha* 33(Suppl. 3), 25–37.
- 540 Bench, J. & Bamford. (1979). *Speech-Hearing Tests and the Spoken Language of*  
541 *Hearing-Impaired Children*. London, Academic Press.
- 542 Byrne, D., Dillon, H., Tran, K.V., Arlinger, S., Wilbraham, K., Cox, R.M., Hagerman,  
543 B., Héту, R., Kei, J., Lui, C.P., Kiessling, J., Kotby, M.N., Nasser, N.H., Kholy,  
544 W.A., Nakanishi, Y., Oyer, H.J., Powell, R., Stephens, D., Meredith, R.,  
545 Sirimanna, T., Tavartkiladze, G., Frolenkov, G.I., Westerman, S., & Ludvigsen,  
546 C. (1994). An international comparison of long-term average speech spectra.  
547 *Journal of the Acoustical Society of America*, 96, 2108-2120.
- 548 Chen, Y., & Wong, L. L. N. (2020). Development of the Mandarin hearing in noise test  
549 for children. *International Journal of Audiology*, 59(9), 707–712.
- 550 China State Language Commission. (2021). *The Report on the Development of Chinese*  
551 *Language*. Beijing, The Commercial Press.
- 552 Ching, T. Y. C., Dillon, H., Leigh, G., & Cupples, L. (2018). Learning from the  
553 Longitudinal Outcomes of Children with Hearing Impairment (LOCHI) study:  
554 summary of 5-year findings and implications. *International Journal of*  
555 *Audiology*, 57(sup2), S105-S111.
- 556 Dillon, H. (1982). A quantitative examination of the sources of speech discrimination  
557 test score variability. *Ear and Hearing* (3), 51–58.
- 558 Dillon H. (1983). The effect of test difficulty on the sensitivity of speech discrimination  
559 tests. *The Journal of the Acoustical Society of America*, 73(1), 336–344.
- 560 Dillon, H. & Ching, T. (1995). What makes a good speech test? in G. Plant & K. E.  
561 Spens (Ed). *Profound Deafness and Speech Communication*. London, Whurr  
562 Publishers: 305–344.

- 563 Etymotic Research. (2005). BKB-SIN Speech-in-Noise Test (Compact Disk), Elk  
564 Grove Village, IL.
- 565 Fu, Q. J., Zhu, M., & Wang, X. (2011). Development and validation of the Mandarin  
566 speech perception test. *The Journal of the Acoustical Society of America*, 129(6),  
567 EL267–EL273.
- 568 Hagerman, B. (1976). Reliability in the determination of speech discrimination.  
569 *Scandinavian Audiology* 5: 219-228.
- 570 Hu, H., Xi, X., Wong, L., Hochmuth, S., Warzybok, A., & Kollmeier, B. (2018).  
571 Construction and evaluation of the Mandarin Chinese matrix (CMNmatrix)  
572 sentence test for the assessment of speech recognition in noise. *International*  
573 *Journal of Audiology*, 57(11), 838–850.
- 574 Keidser, G., Ching, T.Y., Dillon, H., Agung, K., Brew, C., Brewer, S., Fisher, M.J.,  
575 Foster, L., Grant, F., & Storey, L. (2002). The National Acoustic Laboratories  
576 (NAL) CDs of speech and noise for evaluation: normative data and potential  
577 applications. *Australian and New Zealand Journal of Audiology*, 24, 16–35.
- 578 Killion, M. C., Niquette, P. A., Gudmundsen, G. I., Revit, L. J., & Banerjee, S. (2004).  
579 Development of a quick speech-in-noise test for measuring signal-to-noise ratio  
580 loss in normal-hearing and hearing-impaired listeners. *The Journal of the*  
581 *Acoustical Society of America*, 116(4 Pt 1), 2395–2405.
- 582 Leibold, L. J., & Buss, E. (2019). Masked Speech Recognition in School-Age Children.  
583 *Front. Psychol.* 10, 1981.
- 584 Li, J. N., Chen, S., Zhai, L., Han, D. Y., Eshraghi, A. A., Feng, Y., Yang, S. M., & Liu,  
585 X. Z. (2017). The advances in hearing rehabilitation and cochlear implants in  
586 China. *Ear and Hearing*, 38(6), 647–652.
- 587 Litovsky, R. Y., Johnstone, P. M., & Godar, S. P. (2006). Benefits of bilateral cochlear  
588 implants and/or hearing aids in children. *International journal of audiology*, 45  
589 Suppl 1(Suppl 1), S78–S91.
- 590 Liu, S., Han, D.M., Zhang, N., Wu, X., Sheng, Y. L., Mo, L.Y., Yang, X.L., Kong, Y.  
591 (2008). Development of a database of everyday speech material for preschool  
592 Mandarin children (in Chinese). *Journal of Audiology and Speech Pathology*  
593 (China), 16, 121–124
- 594 Moore, B.C. (2003). Speech processing for the hearing-impaired: successes, failures,  
595 and implications for speech mechanisms. *Speech Communication*, 41, 81-91.

- 596 Minimum-Speech-Test-Battery. (2011) The new minimum speech test battery.  
597 <http://auditorypotential.com/MSTB.html>.
- 598 Newman R.S., Morini G, and Ahsan F. (2015) Linguistically-based informational  
599 masking in preschool children. *The Journal of the Acoustical Society of*  
600 *America 138(1), EL93-E198.*
- 601 Nilsson, M., Soli, S. D. & Sullivan, J. (1994). Development of the hearing in noise test  
602 for the measurement of speech reception thresholds in quiet and in noise.  
603 *Journal of the Acoustical Society of America 95, 1085–1099.*
- 604 Nissen, S. L., Harris, R. W., Jennings, L. J., Eggett, D. L., & Buck, H. (2005).  
605 Psychometrically equivalent Mandarin bisyllabic speech discrimination  
606 materials spoken by male and female talkers. *International Journal of*  
607 *Audiology, 44(7), 379–390.*
- 608 Shaughnessy, J. J., Zechmeister, E. B., & Zechmeister, J. S. (2006). *Research methods*  
609 *in psychology*. New York, McGraw Hill.
- 610 Sherbecoe, Robert L. & Studebaker, Gerald A. (2004). Supplementary formulas and  
611 tables for calculating and interconverting speech recognition scores in  
612 transformed arcsine units. *International Journal of Audiology, 43(8), 442 - 448*
- 613 Studebaker, G. A. (1985). A "rationalized" arcsine transform. *Journal of Speech and*  
614 *Hearing Research, 28(3), 455–462.*
- 615 Vaillancourt, V., Laroche, C., Giguère, C., & Soli, S. D. (2008). Establishment of age-  
616 specific normative data for the Canadian French version of the hearing in noise  
617 test for children. *Ear and Hearing, 29(3), 453–466.*
- 618 Wilson, R.H. (2004). Adding speech-in-noise testing to your clinical protocol: Why and  
619 how. *The Hearing Journal, 57, 10.*
- 620 Wilson, R. H., & Carter, A. S. (2001). Relation between slopes of word recognition  
621 psychometric functions and homogeneity of the stimulus materials. *Journal of*  
622 *the American Academy of Audiology, 12(1), 7–14.*
- 623 Wong, L.L, Soli, S.D., Liu, S., Han, N. & Huang, M.W. (2007). Development of the  
624 Mandarin hearing in noise test (MHINT). *Ear and Hearing, 28, 70–74.*
- 625 Wong, L. L. N., Chen, Y., & Leung, K. P. (2019). The Cantonese hearing in noise test  
626 for children. *Trends in Hearing, 23, 1-9.*
- 627 Xi, X., Ching, T. Y., Ji, F., Zhao, Y., Li, J. N., Seymour, J., Hong, M. D., Chen, A. T.,  
628 & Dillon, H. (2012). Development of a corpus of Mandarin sentences in babble

- 629 with homogeneity optimized via psychometric evaluation. *International Journal*  
630 *of Audiology*, 51(5), 399–404.
- 631 Zheng, Y., Soli, D.S, Wang, K., Meng, J., Meng, Z.-L. (2009). Development of the  
632 Mandarin Pediatric Speech Intelligibility (MPSI) test. *International Journal of*  
633 *Audiology*, 48, 718–728.
- 634 Zhu, M.S. & Miao, X.C. (eds.) (1990). *Psycholinguistics*. Shanghai, East China Normal  
635 University Press.

636 **Figure 1.** The cluster of Performance-SNR functions for the 28 sentence lists in babble  
637 noise in adults with normal hearing.

638

639 **Figure 2.** The cluster of Performance-SNR functions for the 27 sentence lists in babble  
640 noise in children with normal hearing.

641

642 **Figure 3.** The scatter diagram of children's performance scores at 3 SNRs with linear  
643 regressions of age and performance at 3 SNRs.

644

645 **Supplemental Material 1.** The SNR levels at which each of the 32 lists were tested for  
646 each subject, with the number of measurements at -3, -5, -7, -9, respectively, and the  
647 total number of measurements at four SNR levels for each list.

648

649 **Supplemental Material 2.** The mean and SD of the 32 lists tested in NH adults at four  
650 SNRs. Lists that were significantly different from others according to ANOVA results  
651 were italicized.

652

653 **Supplemental Material 3.** The SNR levels at which each of the 27 lists were tested for  
654 each subject, with the number of measurements at +1, -2, -5, respectively, and the total  
655 number of measurements at three SNR levels for each list.

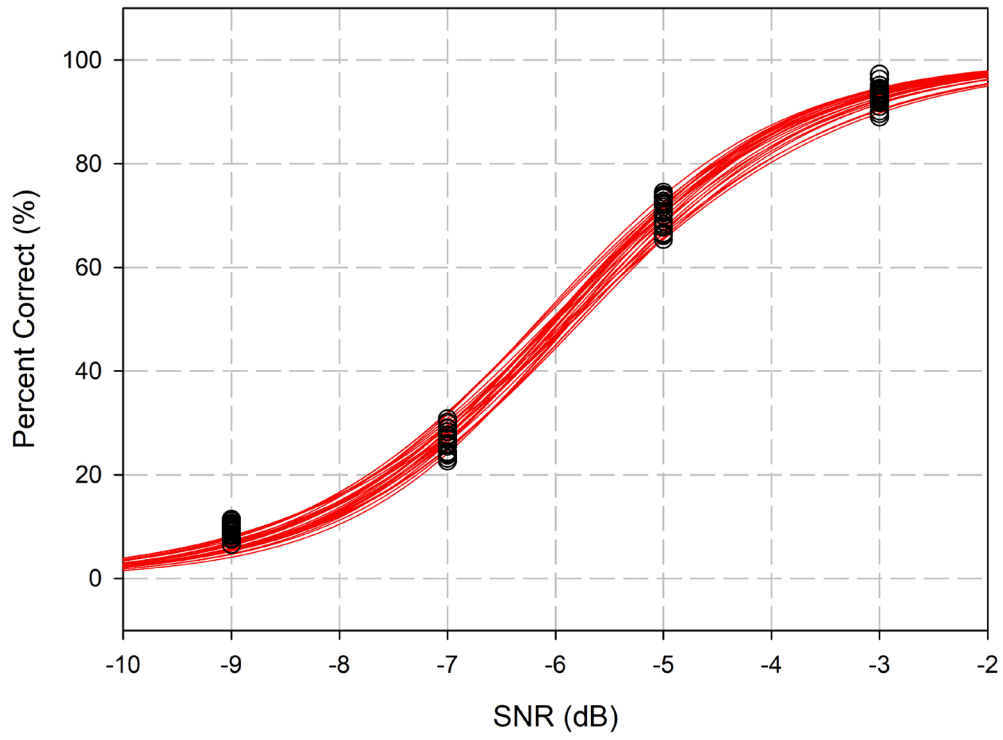
656

657 **Supplemental Material 4.** The mean and SD of the 27 lists tested in NH children at  
658 three SNRs.

659

660 Figure 1.

661



662

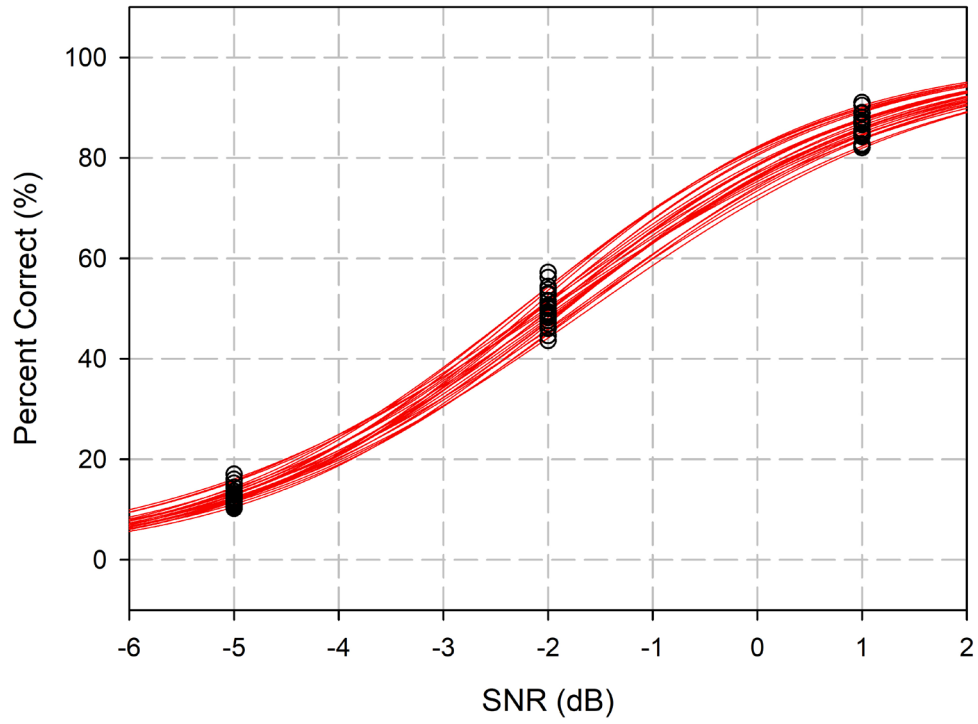
663

664

665

666 Figure 2.

667

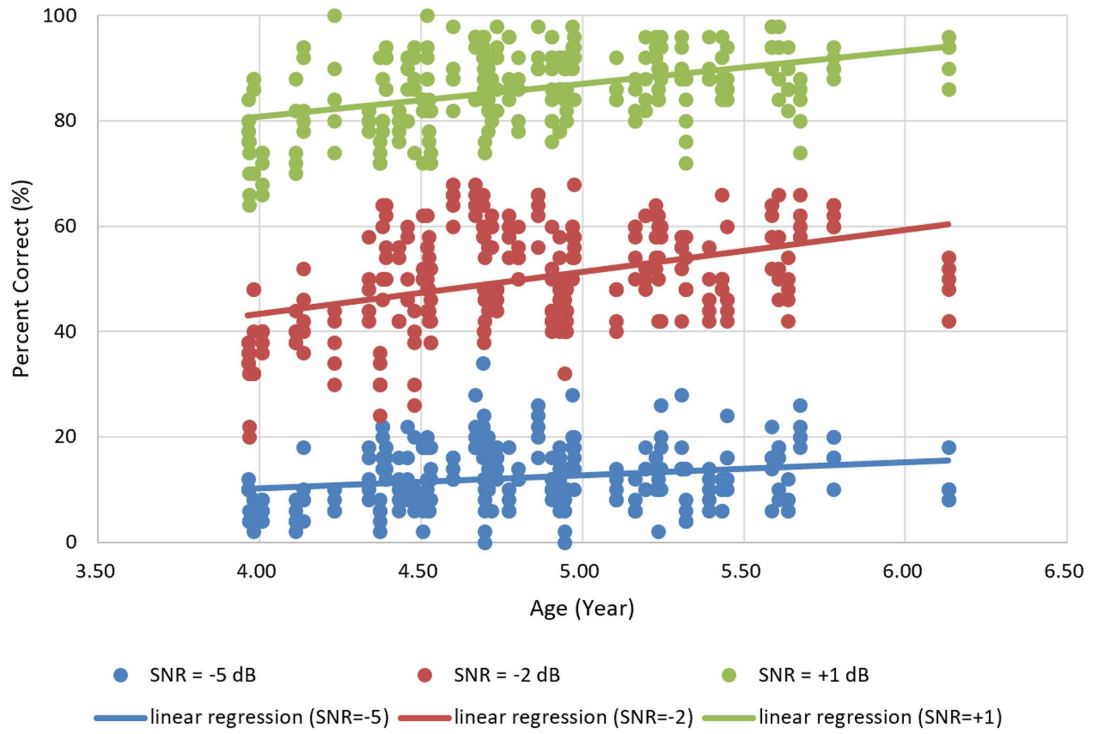


668

669

670 Figure 3.

671



672