## Research Note

# Who is Right? A Word-Identification-in-Noise Test for Young Children Using Minimal Pair Distracters

Samuel Evans[a] and Stuart Rosen[b] iD

[a] Department of Psychology, University of Westminster, London, United Kingdom  [b] Department of Speech, Hearing and Phonetic Sciences, University College London, United Kingdom

ABSTRACT

**Purpose:** Many children have difficulties understanding speech. At present, there are few assessments that test for subtle impairments in speech perception with normative data from U.K. children. We present a new test that evaluates children's ability to identify target words in background noise by choosing between minimal pair alternatives that differ by a single articulatory phonetic feature. This task (a) is tailored to testing young children, but also readily applicable to adults; (b) has minimal memory demands; (c) adapts to the child's ability; and (d) does not require reading or verbal output.
**Method:** We tested 155 children and young adults aged from 5 to 25 years on this new test of single word perception.
**Results:** Speech-in-noise abilities in this particular task develop rapidly through childhood until they reach maturity at around 9 years of age.
**Conclusions:** We make this test freely available and provide associated normative data. We hope that it will be useful to researchers and clinicians in the assessment of speech perception abilities in children who are hard of hearing or have developmental language disorder, dyslexia, or auditory processing disorder.
**Supplemental Material:** https://doi.org/10.23641/asha.17155934

Children with speech, language, and hearing disorders are at a greater risk of poorer literacy (Anthony & Francis, 2005), psychosocial development (Kilpatrick et al., 2019), and long-term prospects (Bryan et al., 2007). Deficits in speech perception, in addition to being a defining feature of hearing impairment and auditory processing disorder (APD; Moore et al., 2013), are associated with a number of developmental disorders, most notably dyslexia (Noordenbos & Serniclaes, 2015) and developmental language disorder (DLD; Ferguson et al., 2011). Developing robust methods to identify individuals with speech perception deficits is a first step toward better characterizing and treating these disorders. At present, there are few tests that assess subtle impairments in speech perception and that have appropriate normative data from U.K. children. Here, we make freely available such a test, which we

envisage will be useful to researchers and clinicians in evaluating the perceptual abilities of young children.

Many children find understanding spoken language difficult. In children who are hard of hearing, these difficulties are obvious and affect perception in both ideal and adverse listening situations. Pure-tone thresholds, although important, provide limited information on functional listening abilities (Houtgast & Festen, 2008), and tests of speech perception in noise provide arguably a more valid assessment of day-to-day listening in children (Leibold et al., 2019). Children with DLDs often exhibit subtle speech perception deficits. However, deficits are not always readily apparent and are sometimes only found in a minority of individuals, or not at all (Messaoud-Galusi et al., 2011). This may reflect a lack of sensitivity of available tests, an absence of a true speech perception deficit, or significant heterogeneity in the individuals assigned to these groups. Only further research will help to uncover which of these explanations is correct. This task is made more difficult by the high comorbidity between developmental reading, language, and auditory

Correspondence to Samuel Evans: S.Evans1@westminster.ac.uk. *Disclosure: The authors have declared that no competing financial or nonfinancial interests existed at the time of publication.*

processing disorders (Bishop et al., 2016; Moore et al., 2013) and the paucity of tools for assessing speech perception in children. A wider range of speech perception tests are required to better characterize the speech perception abilities of children who are hard of hearing and to further our understanding of DLDs.

Successful speech perception requires the integration of multiple covarying acoustic features (Kluender & Alexander, 2010; Lisker, 1977). In natural speech, the multiplicity of available features helps to ensure that perception remains relatively robust to acoustic variation and degradation of the speech signal. Speech sounds that differ on the basis of fewer contrastive features are more highly confusable (Miller & Nicely, 1955). Children with language impairments tend to perform more poorly on tasks in which speech tokens differ minimally from one another such as when categorizing synthetic continua that differ on a single acoustic parameter (Collet et al., 2012; Zoubrinetzky et al., 2016). Deficits in these groups have been shown to be less pronounced in tasks involving natural speech tokens that differ on the basis of multiple acoustic cues (Blomert & Mitterer, 2004; Coady et al., 2005). Speech perception tasks can also be made more challenging by manipulating extrinsic factors, such as the presence of competing noise. Competing sounds generate overlapping patterns of excitation in the auditory periphery that obscure or destroy salient acoustic cues, phenomena referred to as energetic and/or modulation masking (Brungart, 2001; Stone et al., 2011). White noise and steady-state speech spectrum–shaped noise (as used in this study) are expected to interfere with speech perception predominantly through masking of this type. Additional, informational masking effects, those not explained by energetic and modulation masking, are thought to arise at more central, cognitive levels of processing (Shinn-Cunningham, 2008). This form of masking is most often associated with competing speech and is attributable in part to the difficulty of separating out and attending to the correct speech stream.

Speech perception deficits are not always observed in children with DLDs when tested in ideal listening conditions. Performance is often at ceiling, and the addition of competing noise is needed to provide a perceptual stressor that more reliably reveals subtle perceptual deficits (Calcus et al., 2015, 2018; Inoue et al., 2011; Ziegler et al., 2005, 2009). These deficits have been observed in the context of both competing speech (Dole et al., 2012) and competing nonspeech (Ziegler et al., 2005, 2009). Most frequently, deficits have been observed when participants are required to identify and categorize nonword syllables, suggesting a locus of deficit originating at the phonetic and/or phonemic levels (Calcus et al., 2015; Varnet et al., 2016; Ziegler et al., 2005, 2009). Studies have shown weaknesses discriminating specific kinds of phonetic contrasts in children with language impairment (Cornelissen

et al., 1996; Ziegler et al., 2005, 2009). Results from these studies suggest that different language impairments might be associated with deficits in specific phonetic contrasts; for example, children with dyslexia have been shown to have greater difficulty with voicing contrasts, whereas those with DLD have problems with place and manner (Ziegler et al., 2005, 2009). Some studies have also found evidence for generalized deficits, rather than difficulties for specific classes of phonetic contrasts (Calcus et al., 2015).

In typical development, the encoding in the auditory periphery of basic sound features matures early and is thought to be broadly complete by around 6 months of age (Leibold & Buss, 2019). Despite this early maturation, perception in noise abilities continues to mature over a long period. Adultlike perceptual ability does not emerge until 9–10 years of age for speech in steady-state speech-shaped noise (Nishi et al., 2010) and matures even later, around 13–14 years, for speech-in-speech masking (Corbin et al., 2016). This slow development likely reflects the maturation of central auditory and cognitive abilities that relate to sound segregation, dip-listening, selective attention, working memory, and language skills (Leibold & Buss, 2019; Leibold et al., 2019). Young children are easily distracted by additional sound streams, even when the target and masker sounds do not overlap in frequency (Youngdahl et al., 2018). Over time, children learn to deal with distraction and begin to exploit the acoustic distinctions that adults use to improve speech-in-noise performance, such as spatial cues to location (Litovsky, 2005) and differences in pitch and speaker characteristics (Flaherty et al., 2019). Improvements in auditory abilities may also be underpinned by developments in vocabulary and working memory, which have been positively associated with differences in speech-in-noise abilities (McCreery et al., 2017), while noting that these associations have not always been observed (Nittrouer et al., 2013).

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

For the same reason, tests such as the Four Alternative Auditory Feature that require children to read words (Foster & Haggard, 1987) and those using sentences (e.g., Listening in Spatialized Noise–Sentences, Cameron & Dillon, 2007) that place greater demands on auditory working memory and syntactic processing may not always be appropriate. Sentence material may be particularly inappropriate given the evidence that sentence repetition in *quiet* appears to be a good way to diagnose DLD (Conti-Ramsden et al., 2001). Children with language learning impairments such as DLD and dyslexia often have difficulties in reading, syntactic processing, working memory, and vocabulary development (Cowan et al., 2017; Laws et al., 2015; Van Der Lely, 2005). Tests that use single, early acquired words and that require a nonverbal output response allow better assessment of speech perception abilities (especially in young children and those with language learning impairments) as they minimize extraneous syntactic, vocabulary, and working memory demands.

There are relatively few existing U.K. tests of single word perception that have a nonverbal output response. The Consonant Confusion Test (Marriage & Moore, 2003; https://www.chears.co.uk/) is suitable for very young children and requires them to identify a target word from four alternatives presented as pictures. However, in this test, the alternatives differ by multiple phonemes, for example, "cow, owl, house, mouse"; hence, the degree of phonemic discrimination required in this task is relatively broad. The Chear Auditory Perception Test (Marriage & Moore, 2003; https://www.chears.co.uk/) is appropriate for slightly older children and includes contrasts that require a finer level of discrimination. However, the normative data for both these tests are derived from presenting the words at an artificially low volume, used as a way of inducing variation in accuracy (Vickers et al., 2018). This is arguably a less ecologically valid approach, compared to using competing noise to bring accuracy "off ceiling."

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

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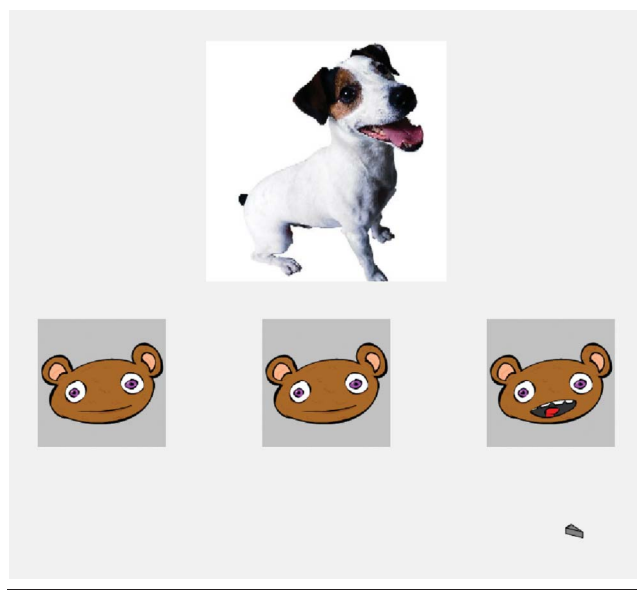poken alternative utterances that are presented against a competing noise. Participants indicate their response nonverbally with a button press. To ensure maximum sensitivity in identifying subtle impairments of speech processing, these alternatives differ by a single articulatory phonetic feature, with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

## Materials and Method

### Test Construction

The WiR? consists of 42 trials, all of a similar form. On each trial, the listener is presented with a picture of a target word on a display screen and hears the same single male speaker produce the name of the target in quiet (see Figure 1). Below, the picture of the target are three cartoon faces that then take turns to speak three utterances. These three utterances are produced by the same single female speaker. Note that the target voice presented in quiet and the voices that participants choose between are from different talkers, intentionally of different sex, so as to prevent participants using an echoic memory trace to perform the task. The voices are presented against a background of steady-state speech spectrum–shaped noise (see details below). Two of the

**Figure 1.** The Who is Right? task. On each trial, the listener sees a picture of a target word and hears the same single male speaker produce the name of the target in quiet. Below, three cartoon faces take turns to speak three utterances presented against a background of steady-state speech spectrum–shaped noise. Two of the utterances are nonword foils differing from the target in a single phonetic feature. The other utterance is the target. Participants select the face that said the "right" word by clicking it with a mouse. A pie chart at the bottom right displays the participant's progress.

utterances are nonword foils differing from the target in its initial consonant in a single feature of voicing, place, or manner (with the two foils always differing in the contrast used). The other utterance is the target. For example, when the target is "bed," the foils are "med" (differing in manner) and "ped" (differing in voicing). The position of the target and two distracter foils are randomized from trial to trial. The listener's task is to identify the face that produced the correct target word by clicking on that face using a mouse. A correct response results in the selected cartoon face smiling, whereas an incorrect response results in the selected face frowning. Every test began with a presentation of 14 familiarization items followed by 28 test items (over which a speech reception threshold [SRT] was calculated), with a random permutation of the items within each phase. All stimuli were presented over headphones at a fixed comfortable level of about 65 dB SPL (measured over the frequency range 100 Hz to 5 kHz).

Target words were monosyllabic words mainly of consonant–vowel–consonant structure (two targets are in consonant–vowel format), which could be presented in an unambiguous pictorial form and whose initial consonant could be altered by a single feature of voicing, manner, or place, to create two nonword foils (see Supplemental Material S1 for full details). All items were early-acquired words, and the test items had a mean age of acquisition of 4.0 years, ranging from 2.9 to 5.6 ($SD$ = 0.67), as measured by Kuperman et al. (2012). For the test trials, the distracter foils comprised 14 manner change items, 21 place change items, and 21 voicing change items, distributed over the 28 test trials (two feature changes per target).

During the test, the signal-to-noise ratio (SNR) was varied adaptively using a two-down/one-up adaptive rule tracking 71% correct (Levitt, 1971), which means that the SNR increases after every error and decreases after two consecutive correct responses. The starting SNR was 20 dB, with a step size of 7 dB, which decreased by 1 dB after every track reversal until it reached 3 dB, at which value it remained for the rest of the test. The SNR was adapted during both the familiarization and test phase. The SRT was defined as the SNR that led to about 71% correct responses, calculated from the mean of the track reversals during the test phase only. Note that lower values indicate better performance, as this indicates that the listener can tolerate poorer SNRs for the desired accuracy. Younger children (under age 9 years) took more time to complete the test, with a median completion time of about 7 min, but everyone older took only about 6 min.

Each test consisted of the same 42 trials (14 familiarization and 28 test items) presented in a different order. The response options on each trial included the target word and the same two unique nonword distracter foils— a stimulus triplet. These stimulus triplets differed greatly in inherent intelligibility, as would be expected by their variety of acoustic, phonetic, and psycholinguistic properties, not to mention the exact choice of foils as being an important

determinant of performance. This is highly undesirable in adaptive testing because it leads to greater variability in the adaptive track. Extensive prior testing on dozens of school-age children (using a combination of adaptive and fixed-SNR testing) allowed the determination of the psychometric functions (relating proportion correct to SNR) for each individual triplet. SRTs for each word were then derived from these functions (through logistic regression), allowing the calculation of a correction factor (the deviation for each triplet from the mean SRT for all triplets) that was applied to the nominal SNR desired during each test (see Supplemental Material S1). This correction factor was used in an additive way to adjust the SNR level up or down for each individual triplet/trial. In this way, performance should be similar for all triplets at the same nominal SNR, which leads to more stable estimates of the SRTs.

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

## Participants

Ethical approval was granted by the University College London Research Ethics Committee. Informed written consent was received from all participants, and their parents, for those aged less than 16 years. None of the children or adults tested had any known speech, hearing, or language impairments, and they were all native British English speakers. These criteria were confirmed by the caregiver during the consent process.

The children and young adults were tested in primary and secondary schools in six separate rounds of testing— referred to as SC ($n$ = 30), GY ($n$ = 17), RL ($n$ = 54), HR ($n$ = 17), HW ($n$ = 18), and CR ($n$ = 19)—and were combined in the analysis. In all instances, testing took place in a quiet room within school, home, or in a quiet, distraction-free public space, for example, a room in a community center. The majority of testing took place in Southern England. Participants for one round of testing (GY) arose from control data from typically developing children as part of a broader study of DLD (Baird et al., 2011; Loucas et al., 2016). Further details concerning the age composition and testing environment for each data set are described in Supplemental Material S2.

There were 155 participants who completed the test (with two exclusions during analysis) and for whom there was complete demographic information (following data exclusions: $M_{age}$ = 11.7 years, ranging from 4.9 to 25.1, $SD$ = 4.6). Gender was well balanced with 63 males and 73 females (54%). There was a mix of genders in all testing rounds. Due to tester error, there were no gender data retained for the CR group, but it was of mixed gender.
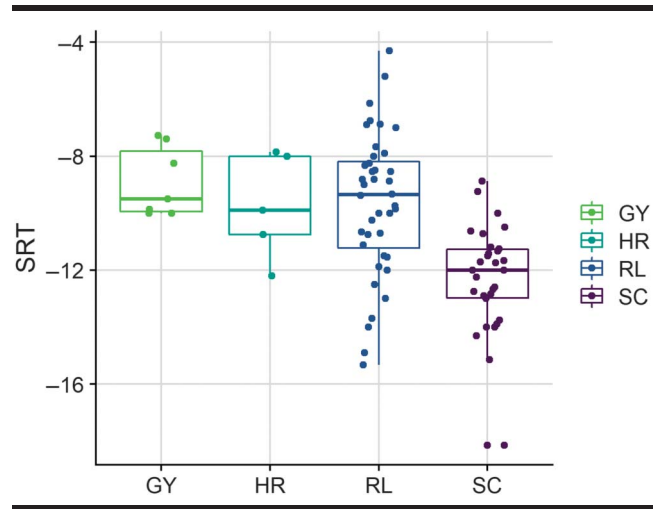
## Results

The mean over the reversals in the test phase of the adaptive track was used to estimate an SRT for each participant. Listeners varied considerably in the total number of reversals that were obtained, from four to 15 ($M$ = 9.6), with 94% of the listeners having seven or more reversals, and no difference, on average, between younger (under 9 years old) and older listeners (within 0.06). There was also no relationship between the number of reversals and age or the SRT. Also of interest is the level of performance observed over the test phase of 28 trials, which should be near the targeted value of 71%. In fact, observed performance levels varied from 61% to 82% ($M$ = 70%), and 95% of listeners had levels within the range of 64%–75%. Again, there was no difference, on average, between younger and older listeners (within 0.5%) and no relationship between performance and age or the SRT. In short, it appears that the adaptive procedure worked equally well across the age range, so any differences in SRT with age likely reflect genuine differences in ability to do the task.

A plot of the obtained data against age showed a strong developmental trend of improving SRTs up to about age 9 or 10 years, leveling off after that point. This also suggested that the SRTs from the SC group (that mainly included older participants) were, on average, better than the other groups for participants of a similar age.

On the basis of the evidence that SRTs did not improve after age 11 years, box plots were made of the SRTs from the four studies for all listeners greater than that age (see Figure 2). A one-way analysis of variance (ANOVA) with a follow-up Tukey's post hoc test confirmed the observation that the mean SRTs were not the same across the four testing groups, $F(3, 78)$ = 9.978, $p$ = $1.22 \times 10^{-5}$. The SRTs for SC were significantly different from RL and GY (both adjusted $ps$ < .003), but SC and HR were not significantly different from each other ($p$ = .086) even though the absolute difference in means was very similar to the other two groups, which did differ. This is likely due to the fact that there are only five older listeners in the HR group.

It is not clear why SRTs were lower in this group, and we assume that this reflects random sampling error. As SC only had participants aged 11.6–16.5 years (in secondary school), it seemed undesirable to leave the SRTs

**Figure 2.** Speech reception thresholds (SRTs) for children aged 11 years and above, illustrating lower SRTs in the SC study. The individual data points are jittered horizontally so as to minimize overlap.



as they were because the overall effect on model fits would not be equal across the age range. Therefore, all SRTs in the SC study were adjusted by the mean difference between the SRTs in that study and the three other studies for children ≥ 11 years old only (by 2.74 dB). A one-way ANOVA confirmed that there was no evidence for differences across the groups after the adjustment, $F(3, 78)$ = 0.256, $p$ = .857.

On the evidence that SRTs change up to about age 9 or 10 years and then asymptote, two different models were used to fit the data. One was a segmented or broken stick regression, in which the model consists of two straight lines that meet at a break point. Two participants were removed from the data set as they contributed a residual with $z$ scores > 3. Once those points were excised, all other $z$ scores were within ± 3. In this fit, a model in which the upper line had a slope = 0 after the break point (implying no change in SRTs after a particular age) was statistically indistinguishable from a model with nonzero slope for the upper segment ($p$ > .4). Also, the broken stick was a much better fit than that provided by a simple linear relationship of SRT with age ($p$ = $3.7 \times 10^{-12}$). The break point was estimated at 9.2 years (95% CI [8.3, 10.2]). Note that, for completeness, the data were also analyzed without the adjustment accounting for the lower SRTs in the SC study and the findings were similar, with a break point at age 10.1 years.

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

$$SRT = b_1 + b_2 * \exp(b_3 * age), \qquad (1)$$

where $b_1$ represents the asymptotic value (i.e., the lowest SRT reached through development), as long as $b_3$ < 0,

Evans & Rosen: Who is Right? Test **5**

which was indeed the case; $b_3$ controls how fast SRTs change over age, and $b_2$ scales the total range of this change. Note the important interaction between $b_2$ and $b_3$ in determining the shape of the curve, whereas $b_1$ is a simple additive term.

The overall fits of the two models were identical, as shown in Figure 3, with a residual standard error of 2.42 on 150 degrees of freedom (as both models have the same number of estimated parameters). We prefer the broken stick model because it gives an unambiguous estimated age for which performance in this task is adultlike. Visualization of the standardized residuals against age for the broken stick regression indicated that variability in measurement of SRT was relatively constant across age after 5 years (see Figure 4).

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

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

*If age ≤ 9.2, Predicted SRT = −1.64 × age + 5.57*
*If age > 9.2, Predicted SRT = −9.6*

Then, a residual is calculated by subtracting the predicted SRT from the actual SRT. This indicates by how many dB a listener is better or worse than an age-matched peer, with negative numbers again indicating better performance. This is then expressed as a *z* score by dividing by

**Figure 4.** The standardized residuals from the broken stick regression.



**Figure 3.** Regression models of speech reception threshold (SRT) with age. The color of the data points indicates which data set they arise from. The two continuous black lines show the predictions of an asymptotic regression (the curved line) and the broken stick regression ("broken" line).
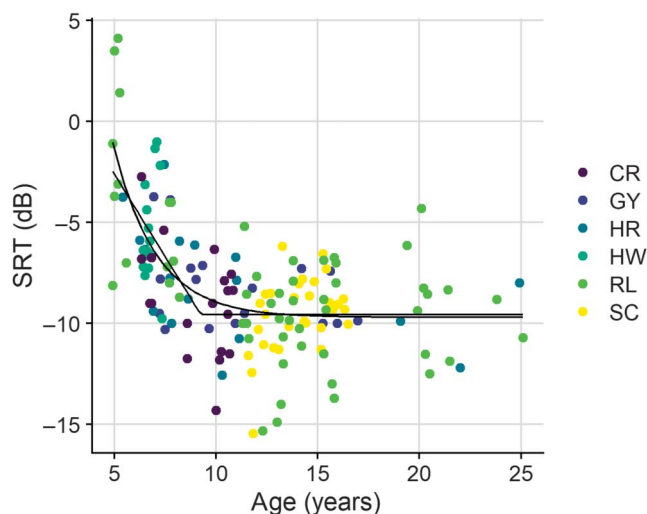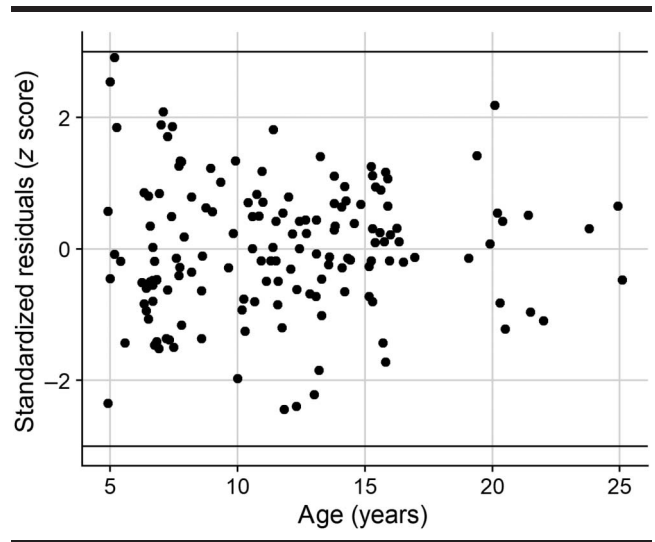


an estimate of the standard deviation of the residuals (2.41). From the *z* score, a percentile can be calculated.

Suppose, for example, that a child aged 6 years obtained an SRT of −0.6 dB. The predicted SRT would be −4.2 dB from the equation above, which means this child is 3.6-dB worse than expected. Dividing through by 2.41 gives *z* = 1.5, which is to say, 1.5 *SD*s worse than typical 6-year-olds. Only about 7% of children of that age would be expected to have an SRT this poor or worse. The test software outputs SRT values in dB, with an option of an extra step to calculate *z* scores based on specifying the listener's age.

## Discussion

We have presented normative data from U.K. children on a test of word identification in noise using minimal pair distracters. A broken stick regression showed that perceptual abilities on this task continued to improve rapidly until the age of around 9 years, before leveling out. We make this task and associated normative data freely available and hope that this test will be of use to researchers and clinicians in the assessment of speech perception abilities of children with language impairments and those that are hard of hearing. In the following sections, we discuss future developments and limitations of the task.

Native language speech sound representations are relatively well developed by 24 months of age but continue to be further refined well into later childhood (Kuhl, 2011). However, the point at which they achieve full maturity is still unknown. Changes are observed until at least 6 years of age (Nittrouer, 2002; Nittrouer & Studdert-Kennedy,

1987), with some studies showing that maturation continues beyond the early teens (Hazan & Barrett, 2000) and into the late teenage years (Davis et al., 2019; McMurray et al., 2018). In the WiR? test, performance rapidly improves until around 9–10 years, before reaching a plateau. This break point is very similar to that obtained in a similar open-response word recognition task in speech spectrum noise in a U.S. sample (Corbin et al., 2016) and is broadly aligned with other studies showing rapid development of speech-in-noise abilities up until the age of 10 years for tasks involving competing energetic/modulation maskers (Hall et al., 2002; Leibold & Buss, 2013; Nishi et al., 2010; Wightman & Kistler, 2005).

The earlier maturation on this task compared with the tasks described above in which maturation continues into the late teenage years (Davis et al., 2019; Hazan & Barrett, 2000; McMurray et al., 2018), may be attributed to important task differences. Our task requires participants to discriminate between canonical articulations with perceptual ambiguity arising from an extrinsic source, the presence of competing noise. By contrast, categorical perception paradigms require participants to categorize ambiguous sounds that are synthesized to be intermediate between canonical articulations. This may require a finer level of phonetic discrimination or place differing demands on decision-making and executive function that give rise to a different developmental trajectory.

The early plateau in energetic masking abilities stands in contrast to the more protracted development associated with informational masking, with adultlike performance on these tasks not achieved until much later, often beyond 13 years of age (Corbin et al., 2016; Hall et al., 2002; Leibold & Buss, 2013). There is also, albeit weak evidence, that SRTs for speech-on-speech masking are a better predictor than equivalent noise masking thresholds for the everyday listening challenges that children who are hard of hearing face (Hillock-Dunn et al., 2015). Such notions may make it seem desirable to implement our task with informational maskers like speech. At present, there is not a speech-on-speech task for children that has normative data from U.K. children. Although it would be possible to construct such a task based on the WiR?, there seems little point to using such carefully constructed stimuli (with the emphasis on the perception of fine phonetic detail), in a version of the task in which higher order abilities like resistance to distraction and auditory scene analysis are important factors. An approach based on simple closed-set targets (e.g., as in Brungart, 2001) might be more appropriate in this instance.

What might be a more promising avenue for these materials, given the different minimal pair contrasts available in WiR?, is to collect normative data on the perception of specific phonetic contrasts. The ability to identify the contrasts that children find most difficult may provide a perspective on the mechanisms that underlie their speech perception weaknesses and allow better targeted interventions for children who are hard of hearing or have DLDs. However, it is likely that such tests would require a fixed SNR, rather than an adaptive approach, with the SNR being fixed at a level appropriate for the listener. In this way, it could be assured that listeners would be not performing near floor or ceiling, but obtain intermediate levels of performance that would allow a sufficient number of errors for meaningful comparisons across contrast types.

The task in its current form also has limitations. At present, we do not have a measure of retest reliability or an understanding of how performance on the test changes with repetitive testing. We hope that retest reliability would be relatively high given the efforts made to calibrate the task through the estimation of an SNR correction factor for each item. Visualization of the standardized residuals of our normative data shows that they are relatively uniformly distributed, with few outliers suggesting that the SRT measure is relatively stable across age. We anticipate that learning in the task would be minimal both within a single test session and across multiple sessions due to the relatively large number of test words and the fact that they are not repeated. Future work addressing retest reliability and learning effects will help to clarify our intuitions. As part of that investigation, it would be useful to know whether it is better to take the first attempt or to average over multiple SRT estimates to attain a truer estimate of speech perception abilities. Indeed, there is some noticeable individual variation in SRT scores (around a 5- to 10-dB range), and greater reliability might be attained by averaging over three measurements (cf. Calandruccio et al., 2020).

Another limitation is that we did not test the pure-tone thresholds for our children and so do not have an objective measure of hearing thresholds for the children in our normative sample. However, all parents reported that their children were without hearing difficulties or speech and language impairments, and we have no reason to think that our sample is unrepresentative of typically developing children. Our full sample (excluding outliers) was 153 participants, a sample size roughly in keeping with or larger than similar tests (Spyridakou et al., 2020; Vance et al., 2009; Vickers et al., 2018). As with most tools of this kind, it would benefit from a larger normative sample and from a broader demographic; factors like social economic status have been shown to influence speech perception ability (Nittrouer, 1996). Our data were collected from only a small number of settings and likely represent a relatively homogenous demographic sample. In the future, normative data from a wider demographic including hard-to-reach populations are necessary, taking into account the additional time and resources that this would entail (Bonevski et al., 2014). As part of this widening

inclusion, it would also be beneficial to consider stratifying by U.K. region to account for differences in regional accent (Adank et al., 2009).

Finally, these normative data apply to quiet listening environments, as might be found in a quiet room within a school or a community clinic. In the future, it would be useful to generate equivalent normative data from children tested in an audiological setting. We hope to address these limitations in the future and allow others to do so, by making this test freely available. We hope that the community will make use of and extend upon our initial work. Only further work will show whether it will be a useful tool in clarifying the speech perception difficulties experienced by listeners with various clinical disorders.

## Acknowledgments

## References

Adank, P., Evans, B. G., Stuart-Smith, J., & Scott, S. K. (2009). Comprehension of familiar and unfamiliar native accents under adverse listening conditions. *Journal of Experimental Psychology: Human Perception and Performance, 35*(2), 520–529. https://doi.org/10.1037/a0013552

Anthony, J., & Francis, D. (2005). Development of phonological awareness skill. *Current Directions in Psychological Science, 14*(5), 255–259. https://doi.org/10.1111/j.0963-7214.2005.00376.x

Baird, G., Slonims, V., Simonoff, E., & Dworzynski, K. (2011). Impairment in non-word repetition: A marker for language impairment or reading impairment? *Developmental Medicine and Child Neurology, 53*(8), 711–716. https://doi.org/10.1111/j.1469-8749.2011.03936.x

Bishop, D. V. M., Snowling, M. J., Thompson, P. A., Greenhalgh, T., Adams, C., Archibald, L., Baird, G., Bauer, A., Bellair, J., Boyle, C., Brownlie, E., Carter, G., Clark, B., Clegg, J., Cohen, N., Conti-Ramsden, G., Dockrell, J., Dunn, J., Ebbels, S., . . . Whitehouse, A. (2016). CATALISE: A multinational and multidisciplinary Delphi consensus study. Identifying language impairments in children. *PLOS ONE, 11*(7), 1–26. https://doi.org/10.1371/journal.pone.0158753

Blomert, L., & Mitterer, H. (2004). The fragile nature of the speech-perception deficit in dyslexia: Natural vs. synthetic speech. *Brain and Language, 89*(1), 21–26. https://doi.org/10.1016/S0093-934X(03)00305-5

Bonevski, B., Randell, M., Paul, C., Chapman, K., Twyman, L., Bryant, J., Brozek, I., & Hughes, C. (2014). Reaching the hard-to-reach: A systematic review of strategies for improving health and medical research with socially disadvantaged groups. *BMC Medical Research Methodology, 14*(1), 1–29. https://doi.org/10.1186/1471-2288-14-42

Brungart, D. S. (2001). Informational and energetic masking effects in the perception of two simultaneous talkers. *The Journal of the Acoustical Society of America, 109*(3), 1101–1109. https://doi.org/10.1121/1.1345696

Bryan, K., Freer, J., & Furlong, C. (2007). Language and communication difficulties in juvenile offenders. *International Journal of Language and Communication Disorders, 42*(5), 505–520. https://doi.org/10.1080/13682820601053977

Byrne, D., Dillon, H., Tran, K., Arlinger, S., Wilbraham, K., Cox, R., Hagerman, B., Hetu, R., Kei, J., Liu, C., Kiessling, J., Notby, M. N., Nasser, N. H. A., El Kholy, W. A. H., Nakanishi, Y., Oyer, H., Powell, R., Stephens, D., Meredith, . . . Ludvigsen, C. (1994). An international comparison of long-term average speech spectra. *The Journal of the Acoustical Society of America, 96*(4), 2108–2120. https://doi.org/10.1121/1.410152

Calandruccio, L., Porter, H. L., Leibold, L. J., & Buss, E. (2020). The clear-speech benefit for school-age children: Speech-in-noise and speech-in-speech recognition. *Journal of Speech, Language, and Hearing Research, 63*(12), 4265–4276. https://doi.org/10.1044/2020_JSLHR-20-00353

Calcus, A., Colin, C., Deltenre, P., & Kolinsky, R. (2015). Informational masking of speech in dyslexic children. *The Journal of the Acoustical Society of America, 137*(6), EL496–EL502. https://doi.org/10.1121/1.4922012

Calcus, A., Hoonhorst, I., Colin, C., Deltenre, P., & Kolinsky, R. (2018). The "Rowdy Classroom Problem in Children with Dyslexia: A Review.". In T. Lachmann & T. Weis (Eds.), *Reading and dyslexia: From basic functions to higher order cognition* (pp. 183–211). Springer International Publishing. https://doi.org/10.1007/978-3-319-90805-2_10

Cameron, S., & Dillon, H. (2007). Development of the Listening in Spatialized Noise–Sentences test (LISN-S). *Ear and Hearing, 28*(2), 196–211. https://doi.org/10.1097/AUD.0b013e318031267f

Coady, J. A., Kluender, K. R., & Evans, J. L. (2005). Categorical perception of speech by children with specific language impairments. *Journal of Speech, Language, and Hearing Research, 48*(4), 944–959. https://doi.org/10.1044/1092-4388(2005/065)

Collet, G., Colin, C., Serniclaes, W., Hoonhorst, I., Markessis, E., Deltenre, P., & Leybaert, J. (2012). Effect of phonological training in French children with SLI: Perspectives on voicing identification, discrimination and categorical perception. *Research in Developmental Disabilities, 33*(6), 1805–1818. https://doi.org/10.1016/j.ridd.2012.05.003

Conti-Ramsden, G., Botting, N., & Faragher, B. (2001). Psycholinguistic markers for specific language impairment (SLI). *Journal of Child Psychology and Psychiatry and Allied Disciplines, 42*(6), 741–748. https://doi.org/10.1111/1469-7610.00770

Corbin, N. E., Bonino, A. Y., Buss, E., & Leibold, L. J. (2016). Development of open-set word recognition in children: Speech-shaped noise and two-talker speech maskers. *Ear and Hearing, 37*(1), 55–63. https://doi.org/10.1097/AUD.0000000000000201

Cornelissen, P. L., Hansen, P. C., Bradley, L., & Stein, J. F. (1996). Analysis of perceptual confusions between nine sets of consonant–vowel sounds in normal and dyslexic adults. *Cognition, 59*(3), 275–306. https://doi.org/10.1016/0010-0277(95)00697-4

Cowan, N., Hogan, T. P., Alt, M., Green, S., Cabbage, K. L., Brinkley, S., & Gray, S. (2017). Short-term Memory in childhood dyslexia: Deficient serial order in multiple modalities. *Dyslexia, 23*(3), 209–233. https://doi.org/10.1002/dys.1557

Davis, M., Evans, S., McCarthy, K., Evans, L., Giannakopoulou, A., & Taylor, J. (2019). Lexical learning shapes the development of speech perception until late adolescence. *PsyArXiv,* 1–52. https://doi.org/10.31234/osf.io/ktsey

Dawes, P., & Bishop, D. V. M. (2007). The SCAN-C in testing for auditory processing disorder in a sample of British children. *International Journal of Audiology, 46*(12), 780–786. https://doi.org/10.1080/14992020701545906

Dole, M., Hoen, M., & Meunier, F. (2012). Speech-in-noise perception deficit in adults with dyslexia: Effects of background type and listening configuration. *Neuropsychologia, 50*(7), 1543–1552. https://doi.org/10.1016/j.neuropsychologia.2012.03.007

Ferguson, M. A., Hall, R. L., Riley, A., & Moore, D. R. (2011). Communication, listening, cognitive and speech perception skills in children with auditory processing disorder (APD) or specific language impairment (SLI). *Journal of Speech, Language, and Hearing Research, 54*(1), 211–227. https://doi.org/10.1044/1092-4388(2010/09-0167)

Flaherty, M. M., Buss, E., & Leibold, L. J. (2019). Developmental effects in children's ability to benefit from f0 differences between target and masker speech. *Ear and Hearing, 40*(4), 927–937. https://journals.lww.com/ear-hearing/Fulltext/2019/07000/Developmental_Effects_in_Children_s_Ability_to.15.aspx

Foster, J. R., & Haggard, M. P. (1987). The Four Alternative Auditory Feature test (FAAF)—Linguistic and psychometric properties of the material with normative data in noise. *British Journal of Audiology, 21*(3), 165–174. https://doi.org/10.3109/03005368709076402

Hall, J. W., Grose, J. H., Buss, E., & Dev, M. B. (2002). Spondee recognition in a two-talker masker and a speech-shaped noise masker in adults and children. *Ear and Hearing, 23*(2), 159–165. https://doi.org/10.1097/00003446-200204000-00008

Hazan, V., & Barrett, S. (2000). The development of phonemic categorization in children aged 6–12. *Journal of Phonetics, 28*(4), 377–396. https://doi.org/10.1006/jpho.2000.0121

Hillock-Dunn, A., Taylor, C., Buss, E., & Leibold, L. J. (2015). Assessing speech perception in children with hearing loss: What conventional clinical tools may miss. *Ear and Hearing, 36*(2), e57–e60. https://journals.lww.com/ear-hearing/Fulltext/2015/03000/Assessing_Speech_Perception_in_Children_With.19.aspx

Houtgast, T., & Festen, J. M. (2008). On the auditory and cognitive functions that may explain an individual's elevation of the speech reception threshold in noise. *International Journal of Audiology, 47*(6), 287–295. https://doi.org/10.1080/14992020802127109

Inoue, T., Higashibara, F., Okazaki, S., & Maekawa, H. (2011). Speech perception in noise deficits in Japanese children with reading difficulties: Effects of presentation rate. *Research in Developmental Disabilities, 32*(6), 2748–2757. https://doi.org/10.1016/j.ridd.2011.05.035

Keith, R. W. (2000). Development and standardization of SCAN-C test for auditory processing disorders in children. *Journal of the American Academy of Audiology, 11*(8), 438–445.

Kilpatrick, T., Leitão, S., & Boyes, M. (2019). Mental health in adolescents with a history of developmental language disorder: The moderating effect of bullying victimisation. *Autism and Developmental Language Impairments, 4*. https://doi.org/10.1177/2396941519893313

Kluender, K. R., & Alexander, J. M. (2010). Perception of speech sounds. In A. I. Basbaum, A. Kaneko, G. M. Shephard, G. Westheimer, T. D. Albright, R. H. Masland, P. Dallos, D. Oertel, S. Firestein, G. K. Beauchamp, C. Bushnell, J. H. Kass, & E. Gardner (Eds.), *The senses: A comprehensive reference*. Elsevier.

Kuhl, P. (2011). Brain mechanisms in early language acquisition. *Neuron, 67*(5), 713–727.

Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for 30,000 English words. *Behavior Research Methods, 44*(4), 978–990. https://doi.org/10.3758/s13428-012-0210-4

Laws, G., Briscoe, J., Ang, S. Y., Brown, H., Hermena, E., & Kapikian, A. (2015). Receptive vocabulary and semantic knowledge in children with SLI and children with Down syndrome. *Child Neuropsychology, 21*(4), 490–508. https://doi.org/10.1080/09297049.2014.917619

Leibold, L. J., & Buss, E. (2013). Children's identification of consonants in a speech-shaped noise or a two-talker masker. *Journal of Speech, Language, and Hearing Research, 56*(4), 1144–1155. https://doi.org/10.1044/1092-4388(2012/12-0011)

Leibold, L. J., & Buss, E. (2019). Masked speech recognition in school-age children. *Frontiers in Psychology, 10*(September). https://doi.org/10.3389/fpsyg.2019.01981

Leibold, L. J., Buss, E., & Calandruccio, L. (2019). Too young for the cocktail party? *Acoustics Today, 15*(1), 37. https://doi.org/10.1121/at.2019.15.1.39

Levitt, H. (1971). Transformed up–down methods in psychoacoustics. *The Journal of the Acoustical Society of America, 49*(2B), 467–477. https://doi.org/10.1121/1.1912375

Lisker, L. (1977). Rapid versus rabid—Catalog of acoustic features that may cue distinction. *The Journal of the Acoustical Society of America, 62,* S77–S78.

Litovsky, R. Y. (2005). Speech intelligibility and spatial release from masking in young children. *The Journal of the Acoustical Society of America, 117*(5), 3091–3099. https://doi.org/10.1121/1.1873913

Loucas, T., Baird, G., Simonoff, E., & Slonims, V. (2016). Phonological processing in children with specific language impairment with and without reading difficulties. *International Journal of Language & Communication Disorders, 51*(5), 581–588. https://doi.org/10.1111/1460-6984.12225

Marriage, J., & Moore, B. (2003). New speech tests reveal benefit of wide-dynamic-range, fast-acting compression for consonant discrimination in children with moderate-to-profound hearing loss. *International Journal of Audiology, 42*(7), 418–425. https://doi.org/10.3109/14992020309080051

McCreery, R. W., Spratford, M., Kirby, B., & Brennan, M. (2017). Individual differences in language and working memory affect children's speech recognition in noise. *International Journal of Audiology, 56*(5), 306–315. https://doi.org/10.1080/14992027.2016.1266703

McMurray, B., Danelz, A., Rigler, H., & Seedorff, M. (2018). Speech categorization develops slowly through adolescence. *Developmental Psychology, 54*(8), 1472–1491. https://doi.org/10.1037/dev0000542

Messaoud-Galusi, S., Hazan, V., & Rosen, S. (2011). Investigating speech perception in children with dyslexia: Is there evidence of a consistent deficit in individuals? *Journal of Speech, Language, and Hearing Research, 54*(6), 1682–1701. https://doi.org/10.1044/1092-4388(2011/09-0261)

Miller, G. A., & Nicely, P. A. (1955). An analysis of perceptual confusions among some English consonant. *The Journal of the Acoustical Society of America, 27,* 338–352. https://doi.org/10.1121/1.1907526

Moore, D. R., Rosen, S., Bamiou, D. E., Campbell, N. G., & Sirimanna, T. (2013). Evolving concepts of developmental auditory processing disorder (APD): A British Society of Audiology APD Special Interest Group "white paper." *International Journal of Audiology, 52*(1), 3–13. https://doi.org/10.3109/14992027.2012.723143

Nishi, K., Lewis, D. E., Hoover, B. M., Choi, S., & Stelmachowicz, P. G. (2010). Children's recognition of American English consonants in noise. *The Journal of the Acoustical Society of America, 127*(5), 3177–3188. https://doi.org/10.1121/1.3377080

Nittrouer, S. (1996). The relation between speech perception and phonemic awareness: Evidence from low-SES children and children with chronic OM. *Journal of Speech and Hearing Research, 39*(5), 1059–1070. https://doi.org/10.1044/jshr.3905.1059

Nittrouer, S. (2002). Learning to perceive speech: How fricative perception changes, and how it stays the same. *The Journal of the Acoustical Society of America, 112*(2), 711–719. https://doi.org/10.1121/1.1496082

Nittrouer, S., Caldwell-Tarr, A., Tarr, E., Lowenstein, J., Rice, C., & Moberly, A. (2013). Improving speech-in-noise recognition for children with hearing loss: Potential effects of language abilities, binaural summation, and head shadow. *International Journal of Audiology, 52*(8), 513–525. https://doi.org/10.3109/14992027.2013.792957

Nittrouer, S., & Studdert-Kennedy, M. (1987). The role of coarticulatory effects in the perception of fricatives by children and adults. *Journal of Speech and Hearing Research, 30*(3), 319–329. https://doi.org/10.1044/jshr.3003.319

Noordenbos, M. W., & Serniclaes, W. (2015). The categorical perception deficit in dyslexia: A meta-analysis. *Scientific Studies of Reading, 19*(5), 340–359. https://doi.org/10.1080/10888438.2015.1052455

Shinn-Cunningham, B. G. (2008). Object-based auditory and visual attention. *Trends in Cognitive Sciences, 12*(5), 182–186. https://doi.org/10.1016/j.tics.2008.02.003

Spyridakou, C., Rosen, S., Dritsakis, G., & Bamiou, D. E. (2020). Adult normative data for the speech in babble (SiB) test. *International Journal of Audiology, 59*(1), 33–38. https://doi.org/10.1080/14992027.2019.1638526

Stone, M. A., Füllgrabe, C., Mackinnon, R. C., & Moore, B. C. J. (2011). The importance for speech intelligibility of random fluctuations in "steady" background noise. *The Journal of the Acoustical Society of America, 130*(5), 2874–2881. https://doi.org/10.1121/1.3641371

Summerfield, Q., Palmer, A. R., Foster, J. R., Marshall, D. H., & Twomey, T. (1994). Clinical evaluation and test–retest reliability of the IHR-McCormick Automated Toy Discrimination Test. *British Journal of Audiology, 28*(3), 165–179. https://doi.org/10.3109/03005369409086564

Vance, M., Rosen, S., & Coleman, M. (2009). Assessing speech perception in young children and relationships with language skills. *International Journal of Audiology, 48*(10), 708–717. https://doi.org/10.1080/14992020902930550

Van Der Lely, H. K. J. (2005). Domain-specific cognitive systems: Insight from grammatical-SLI. *Trends in Cognitive Sciences, 9*(2), 53–59. https://doi.org/10.1016/j.tics.2004.12.002

Varnet, L., Meunier, F., Trollé, G., & Hoen, M. (2016). Direct viewing of dyslexics' compensatory strategies in speech in noise using auditory classification images. *PLOS ONE, 11*(4), e0153781. https://doi.org/10.1371/journal.pone.0153781

Vickers, D. A., Moore, B. C. J., Majeed, A., Stephenson, N., Alferaih, H., Baer, T., & Marriage, J. E. (2018). Closed-set speech discrimination tests for assessing young children. *Ear and Hearing, 39*(1), 32–41. https://doi.org/10.1097/AUD.0000000000000528

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

Youngdahl, C., Healy, E., Yoho, S., Apoux, F., & Rachael, H. (2018). The effect of remote masking on the reception of speech by young school-age children. *Journal of Speech, Language, and Hearing Research, 61*(2), 420–427. https://doi.org/10.1044/2017_JSLHR-H-17-0118

Ziegler, J., Pech-Georgel, C., George, F., Alario, F. X., & Lorenzi, C. (2005). Deficits in speech perception predict language learning impairment. *Proceedings of the National Academy of Sciences of the United States of America, 102*(39), 14110–14115. https://doi.org/10.1073/pnas.0504446102

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

Zoubrinetzky, R., Collet, G., Serniclaes, W., Nguyen-Morel, M. A., & Valdois, S. (2016). Relationships between categorical perception of phonemes, phoneme awareness, and visual attention span in developmental dyslexia. *PLOS ONE, 11*(3), e0151015. https://doi.org/10.1371/journal.pone.0151015