



**Applying Item Response Theory to the Development of  
a Screening Adaptation of the Goldman-Fristoe Test of  
Articulation-2**

Journal:	<i>Journal of Speech, Language, and Hearing Research</i>
Manuscript ID	JSLHR-L-16-0392.R1
Manuscript Type:	Research Note
Date Submitted by the Author:	13-Mar-2017
Complete List of Authors:	Brackenbury, Tim; Bowling Green State University, Communication Sciences and Disorders Zickar, Michael; Bowling Green State University, Psychology Munson, Benjamin; University of Minnesota, Speech-Language-Hearing Sciences Storkel, Holly; University of Kansas, Speech-Language-Hearing: Sciences and Disorders
Keywords:	Speech sound disorders, Assessment, Children

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15 Applying Item Response Theory to the Development of  
16  
17 a Screening Adaptation of the Goldman-Fristoe Test of Articulation-2  
18  
19  
20  
21  
22  
23

24 Tim Brackenbury

25  
26 Michael J. Zickar

27  
28  
29 Bowling Green State University, Ohio  
30  
31

32  
33 Benjamin Munson

34  
35 University of Minnesota, Minneapolis  
36  
37  
38  
39

40  
41 Holly L. Storkel

42  
43 University of Kansas, Lawrence, Kansas  
44  
45  
46  
47

48 Correspondence to: Tim Brackenbury (tbracke@bgsu.edu)  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Abstract

Purpose: Item Response Theory (IRT) is a psychometric approach to measurement that uses latent trait abilities (e.g., speech sound production skills) to model performance on individual items that vary by difficulty and discrimination. An IRT analysis was applied to preschooler's productions of the words on the Goldman-Fristoe Test of Articulation-2 (GFTA-2) to identify candidates for a screening measure of speech sound production skills.

Method: The phoneme accuracies from 154 preschoolers, with speech skills on the GFTA-2 ranging from the 1<sup>st</sup> to above the 90<sup>th</sup> percentile, were analyzed with a two-parameter logistic model.

Results: A total of 108 of the 232 phonemes from stimuli in the sounds-in-words subtest fit the IRT model. These phonemes, and subgroups of the most difficult of these phonemes, correlated significantly with the children's overall percentile scores on the GFTA-2. Regression equations calculated for the five and ten most difficult phonemes predicted overall percentile score at levels commensurate with other screening measures.

Conclusions: These results suggest that speech production accuracy can be screened effectively with a small number of sounds. They motivate further research towards the development of a screening measure of children's speech sound production skills whose stimuli consist of a limited number of difficult phonemes.

1  
2  
3  
4  
5           Screening measures of children’s speech sound production skills typically follow the  
6  
7 protocol of more comprehensive articulation and phonological tests: children are asked to name  
8  
9 pictures of familiar items that elicit the range of consonant phonemes that are typically acquired  
10  
11 during the preschool and early elementary years (e.g., Fluharty, 2001). Pass and fail criteria are  
12  
13 based on phonetic transcriptions of children's productions. Screening tasks typically weigh each  
14  
15 phoneme or word equally, regardless of the ages at which they are typically developed or their  
16  
17 impacts on intelligibility. It is unclear, however, if such a broad-based method is the most  
18  
19 effective or efficient way to meet the purposes of screening, namely, identifying the need for  
20  
21 further evaluation or referral to another professional (ASHA, n.d.). It may be that a shorter word  
22  
23 list, focused on a subset of phonemes in specific word positions, would be better for  
24  
25  
26 distinguishing children with and without potential speech sound disorders (SSD). One challenge  
27  
28 to developing an effective screening instrument for SSD is determining which phonemes in  
29  
30 words best discriminate children with difficulty from those with typical development. The  
31  
32 present study explored this by applying an Item Response Theory (IRT) analysis to 154  
33  
34 children’s productions of the sounds-in-words subtest of one commonly used standardized test of  
35  
36 children's speech production, the Goldman-Fristoe Test of Articulation – 2 (GFTA-2; Goldman  
37  
38 & Fristoe, 2000).

39  
40           IRT is a collection of statistical models that estimate the probability of a person  
41  
42 answering an item correctly based on an estimate of the person’s underlying latent trait as well as  
43  
44 item parameters that relate to features such as discrimination, difficulty, and guessing. By  
45  
46 choosing a particular IRT model, it is possible to better understand how items function, to  
47  
48 develop tailored assessments, and to use a wide variety of psychometric tools (see de Ayala,  
49  
50 2009 for basic information on IRT). Although IRT has been primarily used in educational  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 assessment and psychological research, IRT has been part of research studies in communication  
4  
5 sciences and disorders since the 1980s (Baylor et al., 2011), with a notable increase in its  
6  
7 application over the past 15 years. Recent IRT applications have addressed different aspects of  
8  
9 assessment including, but are not limited to, examinations of performance differences across  
10  
11 populations (Baylor et al., 2013; Baylor et al., 2014; Hula, Doyle, McNeil, & Mikolic, 2006;  
12  
13 Justice, Bowles, & Skibbe, 2006) or multiple forms of the same test (Hoffman, Templin, & Rice,  
14  
15 2012); the precision, weighting, or validity of items within a test (Baylor, Yorkson, Bamer,  
16  
17 Britton, & Amtmann, 2010; Chenault, Berger, Kremer, & Anteunis, 2013; Edmonds & Donovan,  
18  
19 2012; Fergadiotis, Kellough, & Hula, 2015); and the development of a computerized adaptive  
20  
21 version of an existing test (Hula, Kellough, & Fergadiotis, 2015). In addition, and of particular  
22  
23 relevance to the present investigation, IRT has successfully assisted the development of  
24  
25 screening protocols based on existing tests, banks of test items, and previously collected research  
26  
27 data. This work has addressed a wide range of communicative skills, including expressive  
28  
29 language skills of Spanish-speaking preschoolers (Guiberson & Rodriguez, 2014); hearing aid  
30  
31 acceptance, functionality, and use in adults (Chenault, Anteunis, Kremer, & Berger, 2015;  
32  
33 Demorest, Wark, & Erdman 2011; Mokkink, Knol, van Nispen, & Kramer, 2010); participation  
34  
35 across communication contexts by adults with a variety of disorders (Baylor et al., 2013); word  
36  
37 naming in adults with aphasia (del Toro et al., 2011), and vocabulary development in young  
38  
39 children (Makransky, Dale, Havmose, & Bleses, 2016).  
40  
41  
42  
43  
44  
45  
46  
47

48 The current study focused on childhood SSD. Children with SSD present with poor  
49  
50 speech intelligibility as the result of motoric, linguistic, cognitive, sensory, or unspecified issues.  
51  
52 Estimates of their prevalence among preschool and elementary aged children range from 2% to  
53  
54 25% of the general population (Law, Boyle, Harris, Harkness, & Nye, 2000). Clinicians and  
55  
56  
57  
58  
59  
60

1  
2  
3 researchers identify children with SSD through a combination of standardized tests, spontaneous  
4  
5 speech samples, and measures to rule out other causes, such as an oral mechanism examination  
6  
7  
8 to rule out structural anomalies. To date, there have been no published studies examining the use  
9  
10 of IRT to develop an assessment tool for children with SSD. The present study addresses this  
11  
12 need through the following research questions. The first focused on the phonemes identified with  
13  
14 the IRT model. The second and third explored the utility of those phonemes, and subsets of the  
15  
16 phonemes with the greatest difficulty scores, to serve as a screening measure of children's  
17  
18 speech sound production skills.  
19  
20

- 21  
22 1. Which phonemes within the stimuli of the sounds-in-words subtest of the GFTA-2 would fit  
23  
24 within an IRT model?  
25  
26
- 27  
28 2. How well do children's performance on the phonemes in the IRT model, and subsets of those  
29  
30 phonemes, correlate with their percentile score performance on the GFTA-2?  
31  
32
- 33  
34 3. How strongly can children's percentile score performance on the GFTA-2 and identification  
35  
36 as having or not having a speech sound disorder be predicted from their performance on the  
37  
38 phonemes from the IRT model and subsets of those phonemes?  
39

### 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

#### Method

The participants were 154 monolingual boys and girls between 3 and 7 years of age, with and without SSD. This age group was selected because this is an age at which SSD is most likely to be diagnosed, and hence which is subject to a high number of speech and language screening assessments. The participants' data were collected as part of multiple previous research studies conducted by the third and fourth authors (e.g., Munson, Baylis, Krause, & Yim, 2010; Munson & Krause, 2016; Storkel & Hoover, 2010; Storkel, Maekawa & Hoover, 2010). All usable data were included; No potential participants were specifically included or excluded in order to best

1  
2  
3  
4 fit the IRT model. IPA transcriptions were available for each child's productions of the 53 words  
5  
6 on the sounds-in-words subtest of the GFTA-2. Age and percentile scores, however, were only  
7  
8 available for 133 of the participants. These children had a mean age of 57.2 months (4 years, 9  
9  
10 months,  $\pm$  12.73 months) and included 34 3-year-olds, 39 4-year-olds, 44 5-year-olds, 9 6-year-  
11  
12 olds, and 6 7-year-olds (one child's age was not identified).

13  
14  
15 The GFTA-2 (Goldman & Fristoe, 2000) was chosen because it is among the most  
16  
17 widely used standardized tests of children's speech production conducted in used in the 15 years  
18  
19 prior to this study. Its norming sample includes children with and without SSD between the ages  
20  
21 of 2 years, 0 months and 21 years, 11 months. Standard scores on the GFTA-2 are based on  
22  
23 children's performances on the sounds-in-words subtest, in which their productions of target  
24  
25 phonemes within a picture naming task are scored as correct or incorrect. The GFTA-2 percentile  
26  
27 score performances of the 133 participants with complete data sets ranged from 1 to 98, with a  
28  
29 mean of 32.92 ( $\pm$  29.55). The children with SSD included articulatory and phonological issues of  
30  
31 unknown origin, not secondary to other sensory or cognitive issues or diagnoses such as  
32  
33 childhood apraxia of speech. As shown in Figure 1, the percentile score performances of these  
34  
35 children, at each age, reflected the GFTA-2's distribution in which progressively fewer children  
36  
37 scored at lower ends of the percentile score range. The GFTA-2 was administered and scored  
38  
39 using its standard method, in which children are prompted to name pictures. Children who do  
40  
41 not name pictures spontaneously are given a series of progressively greater support until they  
42  
43 produce the target word. Responses are phonetically transcribed.

44  
45  
46 Each of the 154 participants' attempts at the 232 individual phonemes included in the  
47  
48 GFTA-2's sounds-in-words subtest was treated as a separate item, and scored dichotomously as  
49  
50 correct or incorrect. Because phonemes were nested within individual words (e.g., the /s/ in  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 *house* was discrete from the /s/ in *stars*), consonant clusters were categorized by their constituent  
4  
5 phonemes (e.g., *stars* included separate entries for /s/, /t/, /a/, /r/, and /z/). Items that were  
6  
7 answered incorrectly by fewer than 5 participants were eliminated because they did not have  
8  
9 enough variance for analysis. IRTPro 2.0 was used to fit the two-parameter logistic (2PL) model  
10  
11 to all 232 items of the test (Paek & Han, 2012). IRTPro is a software package that estimates IRT  
12  
13 parameters using a variety of possible IRT models. Although there are many models that could  
14  
15 have been chosen for these data, the 2PL was selected because it allows items to vary on both  
16  
17 difficulty and discrimination, two features found to be important in modeling items. In addition,  
18  
19 the 2PL was a reasonable choice given the relatively small sample size. The 2PL model is  
20  
21 represented by the following formula:  
22  
23  
24  
25

$$26 \quad P(u = 1|\theta) = \frac{1}{1 + e^{(-a(\theta-b))}}$$

27  
28 where *a* refers to the item *discrimination* (i.e., the strength of the relation of that item to the  
29  
30 underlying trait), *b* refers to the item *difficulty*, and  $\theta$  refers to the latent trait being measured by  
31  
32 the trait. The 2PL formula uses the item parameters (*a* and *b*) in conjunction with the person  
33  
34 parameter ( $\theta$ ) to predict the probability of answering an item *u* correctly. The parameters, both  
35  
36 item and person, are estimated via IRTPro using maximum likelihood estimation. An iterative  
37  
38 process was carried out to estimate item parameters, eliminating items that did not fit the 2PL  
39  
40 through using the  $\chi^2$  goodness-of-fit statistics estimated by IRTPro. After eliminating poor fit  
41  
42 items, the analysis was re-run, continuing to throw out items until the model fit acceptably well  
43  
44 (i.e., that there were no items that had significant misfit as judged by IRTPro's  $\chi^2$  statistics).  
45  
46  
47  
48  
49  
50  
51

52  
53 Once the IRT analysis was complete, additional statistical analyses were conducted to  
54  
55 identify a) the degree to which the model accounted for the children's overall performances on  
56  
57  
58  
59  
60



1  
2  
3 the GFTA-2 and b) the predictive accuracy of a subset of phonemes from the model to  
4  
5 discriminate children with and without potential SSD.  
6  
7

### 8 **Results**

9  
10 To test whether the data satisfied the requirement of sufficient unidimensionality, a factor  
11 analysis of tetrachoric inter-item correlations was conducted (necessary because the data were  
12 dichotomous) and found that the first factor accounted for 31.5% of the variance in the scale.  
13  
14 This satisfied the requirement that Reckase (1979) identified that the first factor in an exploratory  
15 factor analysis needed to account for at least 25% of the variance to satisfy the unidimensionality  
16 assumption of IRT.  
17  
18  
19  
20  
21  
22  
23

24 The final IRT model consisted of 108 phonemes, which are listed in the Appendix in  
25 order from the highest to lowest difficulty score. They included all of the American English  
26 consonants, except for /h/ and /ʒ/, and the vowels /i, ɪ, ε, æ, ə, ʌ, ə-, aɪ, aʊ/. The consonants, as a  
27 group, occurred in initial and final syllable positions, and as singletons and within clusters. The  
28 phonemes in the 2PL model were from 49 of the 53 words on the GFTA-2 (i.e., all words except  
29 for *ball*, *house*, *ring*, and *thumb*) and included 47 of the 92 phonemes used to determine  
30 percentile scores on the GFTA-2.  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 The analysis began by exploring which area of the underlying trait, commonly denoted  
41 by the Greek letter  $\theta$  in IRT research, provided the most psychometric information. Information  
42 is an IRT-based concept that quantifies the amount of precision provided by the test at varying  
43 levels of  $\theta$ . Traditional measurements of precision, such as standard error of measurement or  
44 reliability, assume that the precision is uniform throughout the range of the trait being measured.  
45  
46 This assumption is likely untrue for many tests given that some tests are designed to be easy, so  
47 that at-risk individuals can be identified, whereas other tests are designed to identify top talent.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Applying IRT to GFTA-2

9

1  
2  
3 In IRT, the *test information function* allows test users to identify at what range of the trait the test  
4 provides precision and at what ranges the test is relatively imprecise. Figure 2 shows the test  
5 information function for the 108-item 2PL model. The most information was provided in the  
6 negative range of the trait continuum, as demonstrated on the left side of the figure in which the  
7 total information values were higher than the standard error, meaning that the test as a whole was  
8 able to provide the most precise measurement at the low ability. There was relatively little  
9 precision at the high end, as shown on the right side of the figure where the standard error  
10 outranked the total information. This suggests that an instrument based on this model would not  
11 be able to distinguish well between children with the very best speech sound production skills  
12 from those at the upper end of the normal range. The greater precision in the negative range is in  
13 line with the goal of using items from the GFTA-2 as a screening test because it emphasizes  
14 differentiating children who are functioning below the normal range from those who are within  
15 the normal range. To increase measurement precision in the positive range, it would be  
16 necessary to write additional items that were high in difficulty and able to discriminate between  
17 average and high ability respondents.

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39 Correlations between the percentile scores on the GFTA-2 from the 133 participants with  
40 complete data sets and their summed accuracy scores of a) the 92 phonemes used to determine  
41 the percentile scores, b) the 108 phonemes in the 2PL model, and c) various subsets of the  
42 phonemes in the model with the greatest difficulty scores are presented in Table 1. Significant  
43 correlations with GFTA-2 percentile scores were found for each of the groups assessed, with  $r^2$ -  
44 values from 0.16 to 0.66. Two sets of multiple regressions were run to determine how well  
45 combinations of the percentile score phonemes and each group of 2PL phonemes contributed to  
46 the children's GFTA-2 percentile scores. In the first set, the percentile phonemes were entered  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 prior to the 2PL phonemes. As shown in Table 1, all but two of the 2PL groups (the 108 and the  
4  
5 2 2PL groups) contributed significantly after the effects of the percentile phonemes were  
6  
7 accounted for ( $p \leq 0.05$ ). In the second set, the 2PL phonemes were entered before the percentile  
8  
9 score phonemes. In this set, the percentile phonemes accounted for significant additional  
10  
11 variance after the 1-, 2-, 3-, and 5-phoneme 2PL groups were entered ( $p < 0.01$ ). However, the  
12  
13 contribution from the percentile phoneme group was not significant after any of the 2PL groups  
14  
15 with 10 or more phonemes were taken into account ( $p > 0.09$ ).  
16  
17  
18  
19

20 The 3-, 5-, and 10-phoneme 2PL groups (see Table 2) were examined as potential  
21  
22 candidates for a screening measure because they were the smallest 2PL groups that accounted for  
23  
24 as much variability in GFTA-2 percentile scores as the 92 phonemes used to determine those  
25  
26 scores ( $r^2 = 0.53, 0.62, 0.67, \text{ and } 0.57$  respectively). This process began by calculating separate  
27  
28 regression equations for each of these groups on the children's GFTA-2 percentile scores. All  
29  
30 three equations were significant at the 0.05 level ( $p < 0.01$ ). The regression equation for the 3-  
31  
32 phoneme group was predicted GFTA-2 percentile score =  $(12.66 * /s/ \text{ in } stars) + (29.76 * /r/ \text{ in}$   
33  
34  $crying) + (12.54 * /θ/ \text{ in } bath) + 7.94$ . The 5-phoneme group's regression equation was predicted  
35  
36 GFTA-2 percentile score =  $(10.38 * /s/ \text{ in } stars) + (21.53 * /r/ \text{ in } crying) + (12.61 * /θ/ \text{ in } bath) +$   
37  
38  $(5.11 * /r/ \text{ in } tree) + (13.61 * /ʃ/ \text{ in } fishing) + 1.81$ . Finally, the 10-phoneme group's regression  
39  
40 equation was predicted GFTA-2 percentile score =  $(6.42 * /s/ \text{ in } stars) + (6.40 * /r/ \text{ in } crying) +$   
41  
42  $(6.05 * /θ/ \text{ in } bath) + (-1.19 * /r/ \text{ in } tree) + (12.21 * /ʃ/ \text{ in } fishing) + (7.35 * /r/ \text{ in } brush) + (9.15 *$   
43  
44  $/ð/ \text{ in } feather) + (9.14 * /ŋ/ \text{ in } monkey) + (16.08 * /r/ \text{ in } rabbit) + (-2.35 * /v/ \text{ in } vacuum) - 0.65$ .  
45  
46  
47  
48  
49  
50

51 Each of these equations was then applied to the 133 participants' productions, yielding  
52  
53 estimated percentile scores. Their utilities as speech screening measures were evaluated by  
54  
55 calculating sensitivity, specificity, and likelihood ratios for cut off points that best approximated  
56  
57  
58  
59  
60

1  
2  
3 1 standard deviation below the mean. Generally, sensitivity and specificity scores  $\geq 80\%$ ,  
4  
5 positive likelihood ratios  $\geq 3$ , and negative likelihood ratios  $\leq 0.3$  are considered preferable  
6  
7 (Dollaghan, 2007). As shown in Table 3, the 3-phoneme regression equation was better at  
8  
9 accurately identifying children performing within the average range than those below (sensitivity  
10  
11 = 62% and specificity = 76%). The 5- and 10-phoneme regression equations outperformed the 3-  
12  
13 phoneme equation, and showed the opposite pattern (with sensitivities at 84% and 88%, and  
14  
15 specificities at 74% and 70%, respectively). The likelihood ratio results also favored the 5- and  
16  
17 10-phoneme regression equations. The positive likelihood ratios were similar for all three  
18  
19 equations, between 2.56 and 3.26, indicating small to moderate probabilities that the children  
20  
21 below the cut off score truly had SSD (Dollaghan, 2007). The negative likelihood ratio of 0.50  
22  
23 for the 3-phoneme equation yielded a mild probability, while the 0.21 and 0.16 results for the 5-  
24  
25 and 10-phoneme equations, respectively, indicating stronger probabilities that children scoring  
26  
27 above the cut off did not have SSD (Dollaghan, 2007). To determine if other phonemes within  
28  
29 the 2PL model would be more accurate, successive blocks of the next 10 and 5 most  
30  
31 discriminating phonemes across the entire model were run using the same process. The results  
32  
33 for all of these calculations were similar to those above, with sensitivity scores consistently 20%  
34  
35 or more lower than specificity scores for the same phonemes.  
36  
37  
38  
39  
40  
41  
42

### 43 Discussion

44  
45  
46 The 2PL model identified 108 phonemes from the stimuli in the sounds-in-words subtest  
47  
48 of the GFTA-2 that significantly discriminated performance for preschool and early elementary  
49  
50 aged children. These included the majority of consonants and vowels in American English, but  
51  
52 did not strongly overlap with the phonemes used by the GFTA-2 to determine percentile scores.  
53  
54  
55 This is not surprising, as the test was “designed to provide a controlled sample of a child’s  
56  
57  
58  
59  
60

1  
2  
3 spontaneous production in words *of the most frequently occurring consonant sounds in Standard*  
4  
5 *American English* [emphasis added]" (Goldman & Fristoe, 2000, pp. 7). In other words, the  
6  
7 phonemes assessed by the GFTA-2 were chosen to represent the wide range of consonant  
8  
9 sounds, not by how well they discriminated performance. In addition, the GFTA-2 scoring  
10  
11 system weighs each phoneme equally, despite the variations in ages at which they are typically  
12  
13 developed or their impacts on intelligibility. These features are similar to other tests based on  
14  
15 classical test theory (e.g., deVellis, 2006).  
16  
17  
18  
19

20 In contrast, the 2PL phonemes and their regression equations align more closely with  
21  
22 item response theory (e.g., Embretson & Reise, 2000) because they include only the phonemes  
23  
24 with the greatest difficulty scores and each phoneme is individually weighted based on its impact  
25  
26 on the predicted score. The phonemes that occurred within the ten most difficult items of the  
27  
28 2PL model were / s, r, θ, ʃ, ð, ŋ, v /. All of these except for / ŋ /, depending on the data source,  
29  
30 are typically later developing phonemes in American English (e.g., Smit, Hand, Freilinger,  
31  
32 Bernthal, & Bird, 1990). The ten most difficult words also included the target phonemes in the  
33  
34 challenging contexts of consonant clusters, medial positions of multisyllabic words, and word  
35  
36 final position. It is likely that these aspects of the target phonemes' difficulty are what  
37  
38 contributed to their potential as screening items, and not simply their inclusion within the GFTA-  
39  
40 2's stimuli. Further, the full 2PL model's inclusion of both easy and difficult phonemes may  
41  
42 explain why it was more precise at discriminating performance at the lower end of the spectrum  
43  
44 than the higher end. A measure that consists of only difficult phonemes may be better at  
45  
46 discriminating performance across the spectrum. Taken together, these results suggest that future  
47  
48 screening measures of children's speech sound productions skills, whether they are or are not  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 developed from existing tests, should consider stimuli that include difficult phonemes in  
4  
5 challenging contexts.  
6  
7

8         The predictive abilities of the 3-, 5-, and 10-phoneme groups were in a positive, but not  
9  
10 overwhelming, direction. As a group, however, they were within the ranges reported for other  
11  
12 assessment of child speech and language disorders. Two systematic reviews of screening  
13  
14 measures of preschooler's speech and language skills (Law, Boyle, Harris, Harkness, & Nye,  
15  
16 2000; Nelson, Nygern, Walker, & Panoscha, 2006), for example, revealed sensitivity ranges  
17  
18 from 17 – 100%, and specificity ranges from 14 – 100%. It is noted, however, that approximately  
19  
20 half of these screening measures identified fell below the suggested 80% lower limits for  
21  
22 sensitivity or specificity (Dollaghan, 2007). Of the three groups assessed in this study, the one  
23  
24 with 3 phonemes appears to be the weakest, due to its poorer sensitivity, specificity, and  
25  
26 likelihood values. The results for the 5- and 10-phoneme groups were both better and fairly  
27  
28 similar to each other. Caution is advised before directly applying the results of this study to  
29  
30 clinical or research settings. Because the regression equations were calculated from a subset of  
31  
32 children used to develop the 2PL model, for example, it is currently unclear how well these  
33  
34 results will generalize to other children. In addition, the concurrent validity of the 5- and 10-  
35  
36 phoneme groups with other standardized measures of speech sound production should be  
37  
38 evaluated. As a result, direct applications of the phonemes and words within the 2PL model to  
39  
40 speech screening are not recommended without additional research.  
41  
42  
43  
44  
45  
46  
47

48         Although successful, the 2PL model was relatively simple, due to its dichotomous  
49  
50 scoring of item responses. With larger data sets, more flexible models could be used to fit these  
51  
52 data, including the 3PL model that allows for guessing, and polytomous IRT models that would  
53  
54 allow graded responses to be scored (see Zickar, 2002). The latter might be useful in  
55  
56  
57  
58  
59  
60

1  
2  
3 determining whether a scoring system that addresses the specific type of errors (such as  
4  
5 phonological process or distinctive feature differences) would help improve the measurement. In  
6  
7 addition, larger sample sizes would allow for us to estimate these more complex models as well  
8  
9 as model some of the easy items that few children answered incorrectly. Additional areas for  
10  
11 future exploration on this topic include comparing the results of similar IRT analyses on other  
12  
13 measures of speech sound production and examining if and how the IRT results may vary across  
14  
15 age groups.  
16  
17  
18  
19  
20  
21

22 Acknowledgments: The data analyzed in this article were collected in studies that were funded  
23  
24 by NIH grant R03 DC005702 to Benjamin Munson and grant R03 DC006545 to Holly L.  
25  
26 Storkel.  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## References

- American Speech-Language-Hearing Association (n.d.). Speech sound disorders – Articulation and phonology.  
<http://www.asha.org/PRPSpecificTopic.aspx?folderid=8589935321&section=Assessment>
- Baylor, C., McAuliffe, M. J., Hughes, L. E., Yorkston, K., Anderson, T., Kim, J., & Amtmann, D. (2014). A differential item functioning (DIF) analysis of the communicative participation item bank (CPIB): Comparing individuals with Parkinson's disease from the United States and New Zealand. *Journal of Speech, Language, and Hearing Research*, *57*, 90-95.
- Baylor, C., Hula, W., Donovan, N. J., Doyle, P. J., Kendall, D., & Yorkson, K. (2011). An introduction to item response theory and Rasch models for speech-language pathologists. *American Journal of Speech Language Pathology*, *20*, 243-259.
- Baylor, C., Yorkston, K., Bamer, A., Britton, D., & Amtmann, D. (2010). Variables associated with communicative participation in people with Multiple Sclerosis: A regression analysis. *American Journal of Speech-Language Pathology*, *19*, 143-153.
- Baylor, C., Yorkston, K., Eadie, T., Kim, J., Chung, H., & Amtmann, D. (2013). The communicative participation item bank (CPIB): Item bank calibration and development of a disorder-generic short form. *Journal of Speech, Language, and Hearing Research*, *56*, 1190-1208.
- Chenault, M., Anteunis, L., Kremer, B., & Berger, M. (2015). An investigation of measurement equivalence in hearing response scales: Refinement of a questionnaire for use in hearing screening. *American Journal of Audiology*, *24*, 188-203.



1 Applying IRT to GFTA-2

16

2  
3 Chenault, M., Berger, M., Kremer, B., & Anteunis, L. (2013). Quantification of experienced  
4 hearing problems with item response theory. *American Journal of Audiology*, 22, 252-  
5  
6 262.  
7  
8

9  
10 de Ayala, R.J. (2009). *The Theory and Practice of Item Response Theory*. New York: Guilford  
11  
12 Press.  
13

14  
15 del Toro, C. M., Bislick, L. P., Comer, M., Velozo, C., Romero, S., Gonzalex Rothi, L. J., &  
16  
17 Kendall, D. L. (2011). *Journal of Speech, Language, and Hearing Research*, 54, 1089-  
18  
19 1100.  
20  
21

22 Demorest, M. E., Wark, D. J., & Erdman, S. A. (2011). Development of the screening test for  
23  
24 hearing problems. *American Journal of Audiology*, 20, 100-110.  
25  
26

27 deVellis, R. F. (2006). Classical test theory. *Medical Care*, 44, S50-S59.  
28

29 Dollaghan, C. A. (2007). *The handbook for evidence-based practice in communication disorders*.  
30  
31 Baltimore, MD: Brookes.  
32  
33

34 Edmonds, L. A. & Donovan, N. J. (2012). Item-level psychometrics and predictors of  
35  
36 performance for Spanish/English bilingual speakers on an object and action naming  
37  
38 battery. *Journal of Speech, Language, and Hearing Research*, 55, 359-381.  
39  
40

41 Embretson, S. E. & Reise, S. P. (2000). *Item Response Theory for Psychologists*. Mahwah, NJ:  
42  
43 Lawrence Erlbaum Associates.  
44  
45

46 Fergadiotis, G., Kellough, S., & Hula, W. D. (2015). Item response theory modeling of the  
47  
48 Philadelphia Naming Test. *Journal of Speech, Language, and Hearing Research*, 58,  
49  
50 865-877.  
51  
52

53 Fluharty, N. B. (2001). *Fluharty Preschool Speech and Language Screening Test*. Austin, TX:  
54  
55 Pro-Ed.  
56  
57  
58  
59  
60

Applying IRT to GFTA-2

17

1  
2  
3 Goldman, R. & Fristoe, M. (2000). *Goldman-Fristoe Test of Articulation – Second Edition*. San  
4  
5 Antonio, TX: Pearson.

6  
7  
8 Guiberson, M. & Rodriguez, B. L. (2014). Rasch analysis of a Spanish language-screening  
9  
10 parent survey. *Research in Developmental Disabilities, 35*, 646–656.

11  
12 Hoffman, L., Templin, J., & Rice, M. L. (2012). Linking outcomes from Peabody Picture  
13  
14 Vocabulary Test forms using item response models. *Journal of Speech, Language, and*  
15  
16 *Hearing Research, 55*, 754-763.

17  
18  
19  
20 Hula, W., Doyle, P. J., McNeil, M. R., & Mikolic, J. M. (2006). Rasch modeling of Revised  
21  
22 Token Test Performance: Validity and sensitivity to change. *Journal of Speech,*  
23  
24 *Language, and Hearing, Research, 49*, 27-46.

25  
26  
27 Hula, W. D., Kellough, S., & Fergadiotis, G. (2015). Development and simulation testing of a  
28  
29 computerized adaptive version of the Philadelphia Naming Test. *Journal of Speech,*  
30  
31 *Language, and Hearing Research, 58*, 878-890.

32  
33  
34 Justice, L. M., Bowes, R. P., Skibbe, L. E. (2006). Measuring preschool attainment of print-  
35  
36 concept knowledge: A study of typical and at-risk 3- to 5-year-old children using item  
37  
38 response theory. *Language, Speech, and Hearing Services in Schools, 37*, 224-235.

39  
40  
41 Law, J., Boyle, J., Harris, F., Harkness, A. & Nye, C. (2007). The feasibility of universal  
42  
43 screening for primary speech and language delay: Findings from a systematic review of  
44  
45 the literature. *Developmental Medicine and Child Neurology, 42*, 190-200.

46  
47  
48 Makransky, G., Dale, P. S., Havmose, P., & Bleses, D. (2016). An item response theory-based,  
49  
50 computerized adaptive testing version of the MacArthur–Bates Communicative  
51  
52 Development Inventory: Words & Sentences (CDI:WS). *Journal of Speech, Language,*  
53  
54 *and Hearing Research, 59*, 281-289.  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Mokkink, L. B., Knol, D. L., van Nispen, R. M. A., & Kramer, S. E. (2010). Improving the  
4  
5 quality and applicability of the Dutch Scales of the Communication Profile for the  
6  
7 hearing impaired using item response theory. *Journal of Speech, Language, and Hearing*  
8  
9 *Research, 53*, 556-571.  
10  
11
- 12 Munson, B., Baylis, A. L., Krause, M. O. & Yim, D. (2010). Representation and access in  
13  
14 phonological impairment. In, Fougeron, C., Kühnert, B., D'Imperio, M., & Vallée (Eds.),  
15  
16 *Laboratory phonology 10*. Berlin: Walter de Gruyter.  
17  
18
- 19 Munson, B., & Krause, M.O.P. (2016). Phonological Encoding in Speech Sound Disorder:  
20  
21 Evidence from a Cross-Modal Priming Experiment. *International Journal of Language*  
22  
23 *and Communication Disorders*, doi: 10.1111/1460-6984.12271. [Epub ahead of print]  
24  
25  
26
- 27 Nelson, H. D., Nygern, P., Walker, M., & Panoscha, R. (2006). Screening for speech and  
28  
29 language delay in preschool children: Systematic evidence review for the US preventive  
30  
31 services task force. *Pediatrics, 117*, e298-e319.  
32  
33
- 34 Paek, I., & Han, K.T. (2012). IRTPRO 2.1 for windows (item response theory for patient-  
35  
36 reported outcomes). *Applied Psychological Measurement, 37*, 242-252.  
37  
38
- 39 Reckase, M. D. (1979). Unifactor latent trait models applied to multifactor tests: Results and  
40  
41 implications. *Journal of Educational Statistics, 4*, 207-230.  
42  
43
- 44 Revelle, W. (2016). *Procedures for psychological, psychometric, and personality research*.  
45  
46 Published online at <http://personality-project.org/r/psych> (Accessed 2/16/2017).  
47  
48
- 49 Smit, A. B., Hand, L., Freilinger, J., Bernthal, J., & Bird, A. (1990). The Iowa Articulation  
50  
51 Norms Project and its Nebraska replication. *Journal of Speech and Hearing Disorders,*  
52  
53 *55*, 779-798.  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Storkel, H. L. & Hoover, J. R. (2010). Word learning by children with phonological delays:

4  
5 Differentiating effects of phonotactic probability and neighborhood density. *Journal of*  
6  
7  
8 *Communication Disorders*, 43, 105-119.

9  
10 Storkel, H. L., Maekawa, J., & Hoover, J. R. (2010). Differentiating the effects of phonotactic  
11  
12 probability and neighborhood density on vocabulary comprehension and production: A  
13  
14 comparison of preschool children with versus without phonological delays. *Journal of*  
15  
16  
17 *Speech, Language, and Hearing Research*, 53, 933-949.

18  
19 Zickar, M. J. (2002). Modeling data with polytomous item response theory. In F. Drasgow & N.  
20  
21 Schmitt (Eds.), *Measuring and analyzing behavior in organizations: Advances in*  
22  
23  
24 *measurement and data analysis* (pp. 123-155). San Francisco, CA: Jossey-Bass.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure 1. Histogram depicting GFTA-2 performance by age and standard deviation score.  
Figure 2. Test Information Function for the 108-item 2PL model.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Applying IRT to GFTA-2

1

Table 1. Multiple stepwise regressions for two sets of predictors, percentile score phonemes and sets of phonemes from the two-parameter logistic model (2PL), on 133 participants' percentile scores on the GFTA-2.

Comparison phoneme group	2PL phonemes entered before percentile score phonemes				Percentile score phonemes entered before 2PL phonemes			
	Run	r	r <sup>2</sup>	p	Run	r	r <sup>2</sup>	p
108 2PL phonemes	1	0.70	0.49	< 0.01	1	0.71	0.50	< 0.01
	2	0.71	0.50	< 0.01	2			0.75
20 most difficult 2PL phonemes	1	0.76	0.57	< 0.01	1	0.71	0.50	< 0.01
	2			0.66	2	0.76	0.57	< 0.01
15 most difficult 2PL phonemes	1	0.76	0.58	< 0.01	1	0.71	0.50	< 0.01
	2			0.79	2	0.76	0.58	< 0.01
10 most difficult 2PL phonemes	1	0.76	0.57	< 0.01	1	0.71	0.50	< 0.01
	2			0.22	2	0.76	0.58	< 0.01
5 most difficult 2PL phonemes	1	0.73	0.53	< 0.01	1	0.71	0.50	< 0.01
	2	0.74	0.55	< 0.01	2	0.74	0.55	< 0.01
3 most difficult 2PL phonemes	1	0.68	0.46	< 0.01	1	0.71	0.50	< 0.01
	2	0.73	0.53	< 0.01	2	0.73	0.53	< 0.01
2 most difficult 2PL phonemes	1	0.66	0.43	< 0.01	1	0.71	0.50	< 0.01
	2	0.72	0.52	< 0.01	2	0.72	0.52	< 0.01
1 most difficult 2PL phoneme	1	0.42	0.18	< 0.01	1	0.71	0.50	< 0.01
	2	0.71	0.51	< 0.01	2			0.13

## Applying IRT to GFTA-2

Table 2. Predicted GFTA-2\* regression equations for the 3-, 5-, and 10-phoneme Groups.

3-phoneme group		5-phoneme group		10-phoneme group	
/s/ in <i>stars</i>	* 12.66	/s/ in <i>stars</i>	* 10.38	/s/ in <i>stars</i>	* 6.42
/r/ in <i>crying</i>	* 29.76	/r/ in <i>crying</i>	* 21.53	/r/ in <i>crying</i>	* 6.40
/θ/ in <i>bath</i>	* 12.54	/θ/ in <i>bath</i>	* 12.61	/θ/ in <i>bath</i>	* 6.05
+	7.94	/r/ in <i>tree</i>	* 5.11	/r/ in <i>tree</i>	* -1.19
		/ʃ/ in <i>fishing</i>	* 13.61	/ʃ/ in <i>fishing</i>	* 12.21
		+	1.81	/r/ in <i>brush</i>	* 7.35
				/ð/ in <i>feather</i>	* 9.15
				/ŋ/ in <i>monkey</i>	* 9.14
				/r/ in <i>rabbit</i>	* 16.08
				/v/ in <i>vacuum</i>	* -2.35
				+	-0.65

\* Goldman-Fristoe Test of Articulation-2

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Applying IRT to GFTA-2

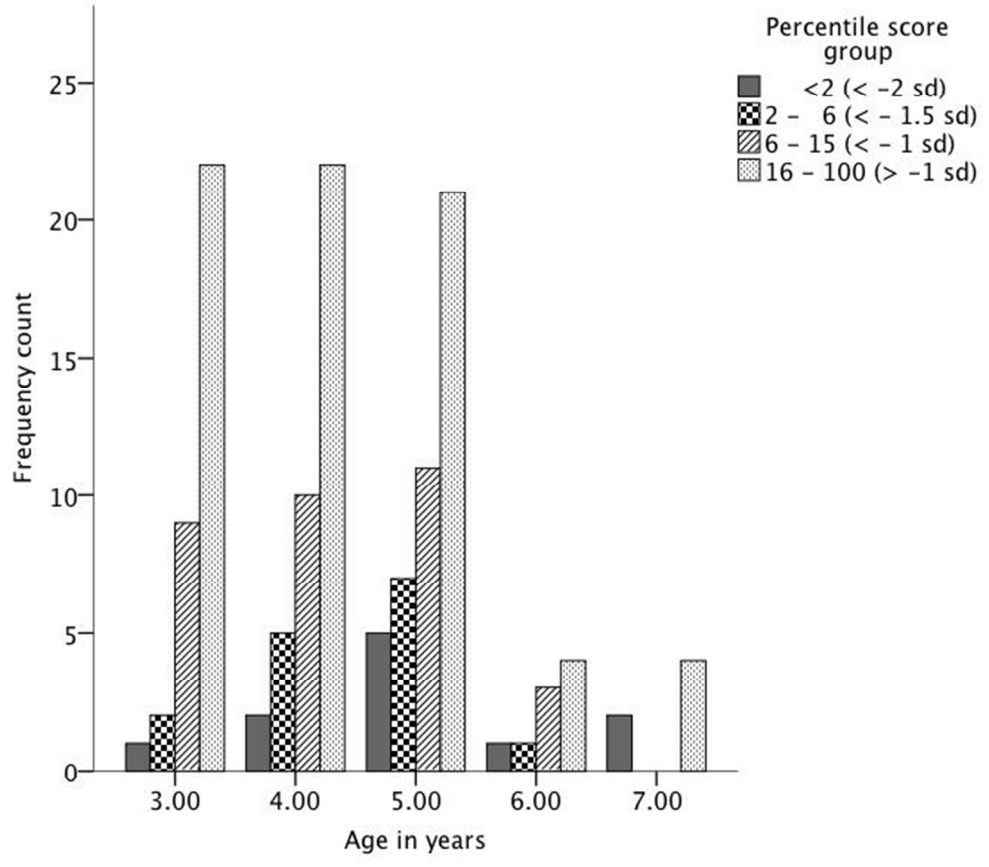
Table 3. Sensitivity, specificity, and likelihood ratio calculations based on three different regression equations developed from the two-parameter logistic model, at a cut off of the 16<sup>th</sup> percentile on the GFTA-2.

Regression equation	Regression cut off score	Sensitivity	Specificity	Positive likelihood ratio	Negative likelihood ratio
3-phoneme	20	62%	76%	2.56	0.50
5-phoneme	17	84%	74%	3.26	0.21
10-phoneme	17	88%	70%	2.98	0.16

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



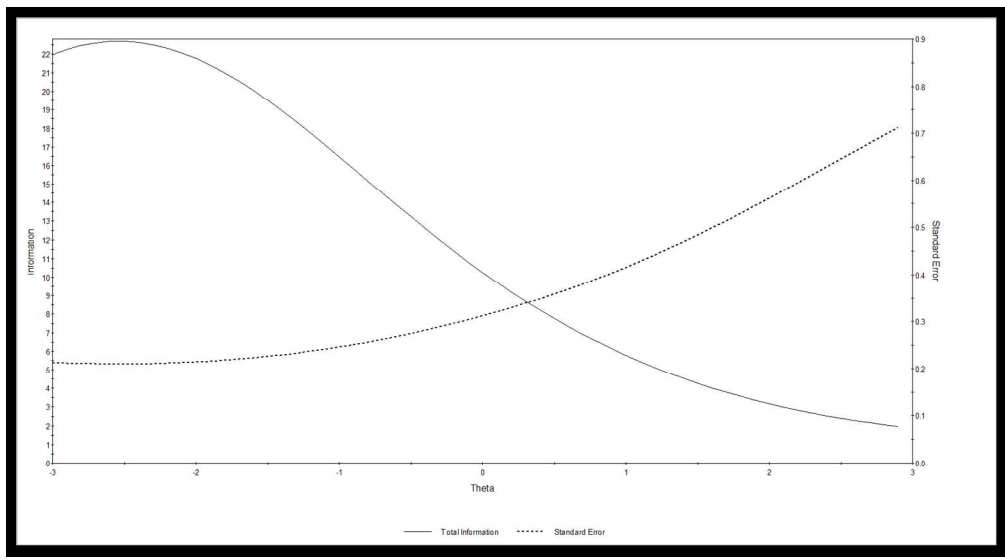
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1

57x50mm (300 x 300 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



2

*Appendix*

The 108 phonemes from the 2PL model, ranked from highest to lowest difficulty score.

	Difficulty		Difficulty
Phoneme	score	Phoneme	score
/s/ in stars*	1.21	/g/ in girl*	0.16
/r/ in crying*	1.15	/ŋ/ in finger*	0.16
/θ/ in bath*	1.15	/l/ in yellow	0.16
/r/ in tree*	1.03	/ə-/ in finger	0.14
/f/ in fishing*	0.92	/l/ in telephone	0.12
/r/ in brush*	0.61	/r/ in orange	0.11
/ð/ in feather*	0.54	/l/ in glasses*	0.11
/ŋ/ in monkey	0.48	/k/ in clown*	0.1
/r/ in rabbit*	0.47	/l/ in shovel	0.1
/v/ in vacuum*	0.43	/ə-/ in scissors	0.08
/r/ in green*	0.41	first /z/ in scissors*	0.08
/r/ in frog	0.4	/k/ in cup*	0.06
/z/ in zipper*	0.4	/k/ in car	0.06
/g/ in wagon*	0.4	/l/ in plane*	0.06
/k/ in crying*	0.37	/z/ in pajamas	0.05
/f/ in shovel*	0.34	/r/ in carrot*	0.05
/ð/ in this*	0.34	/ʌ/ in banana	0.04
/θ/ in bathtub*	0.33	/ə-/ in feather	0.04
/z/ in glasses	0.32	/dʒ/ in orange*	0.03
/s/ in swimming*	0.31	/r/ in stars	0.03
/l/ in watches	0.27	/j/ in vacuum	0.02
/dʒ/ in pajamas*	0.26	/g/ in finger	0
/f/ in five	0.24	/dʒ/ in jumping*	-0.02
/l/ in finger	0.24	/aɪ/ in crying	-0.04
/l/ in balloons*	0.22	/b/ in banana	-0.04
/l/ in flowers*	0.2	/v/ in shovel*	-0.05
first /k/ in quack*	0.19	/f/ in finger	-0.05
/f/ in fishing*	0.19	/j/ in yellow*	-0.06
/ʌ/ in pajamas	0.19	/w/ in swimming*	-0.07
/ə-/ in girl	0.18	/ŋ/ in fishing	-0.08
/p/ in pajamas	0.17	/ŋ/ in jumping	-0.12

*Appendix continued.*

Phoneme	Difficulty score	Phoneme	Difficulty score
/k/ in duck*	-0.12	/t/ in bathtub*	-0.76
/n/ in pencils	-0.13	/æ/ in glasses	-0.85
/tʃ/ in watch*	-0.13	/u/ in vacuum	-0.88
/f/ in feather	-0.15	/d/ in window*	-0.88
/k/ in vacuum	-0.16	/aʊ/ in flowers	-0.91
/l/ in lamp*	-0.17	/b/ in rabbit*	-1.06
/f/ in knife*	-0.21	/u/ in spoon	-1.07
/ʌ/ in balloons	-0.24	/ɪ/ fishing	-1.14
/w/ in watches	-0.33	/aɪ/ in slide	-1.16
/ə/ in wagon	-0.34	/æ/ in bathtub	-1.2
/n/ in orange	-0.36	/ɪ/ in window	-1.22
/ɛ/ in feather	-0.41	/æ/ in lamp	-1.26
/ɛ/ in pencils	-0.41	/n/ balloons	-1.3
/ɪ/ in orange	-0.41	/æ/ in bath	-1.56
/t/ in telephone*	-0.41	/æ/ in rabbit	-1.65
/d/ in drum*	-0.43	/ʌ/ in bathtub	-1.8
/ɪ/ in chair	-0.44	/ʌ/ in shovel	-1.88
/ə/ in pajamas	-0.46	/i/ in green	-1.98
/n/ in telephone	-0.48	/n/ in wagon	-2.46
/ə/ in pencils	-0.55	/b/ in blue*	-3.16
/ɪ/ in scissors	-0.59	/t/ stars*	-3.56
/m/ in vacuum	-0.6	/n/ in knife*	-3.91
/n/ in window	-0.7	/ə/ in banana	-4.5

\* Phonemes included in both the 2PL model and the percentile scoring for the GFTA-2