

1 Running Header: The Houston Questionnaire

2

3 **Can Bilingual Children Self-Report their Bilingual Experience and Proficiency? The**
4 **Houston Questionnaire**

5

6 Anny Castilla-Earls

7 University of Houston

8 Department of Communication Disorders and Sciences

9 Houston, TX

10

11 Juliana Ronderos

12 University of Houston

13 Department of Communication Disorders and Sciences

14 Houston, TX

15

16 Lisa Fitton

17 University of South Carolina

18 Department of Communication Sciences and Disorders

19 Columbia, SC

20

21 Journal-formatted version of this manuscript: https://doi.org/10.1044/2022_JSLHR-21-00675

22 Conflict of Interest: There are no conflicts of interest.

23 Corresponding Author: Anny Castilla-Earls, University of Houston, Melcher Life Sciences

24 3871 Holman St. Room M242; Houston, Tx, 77204. Phone: (713) 743-0488. Email:

25 annycastilla@uh.edu.

26 Funding: Research reported in this publication was supported by the National Institute on

27 Deafness and Other Communication Disorders of the National Institutes of Health under Award

28 Number K23DC015835 granted to Anny Castilla-Earls. The content is solely the responsibility

29 of the authors and does not necessarily represent the official views of the National Institutes of

30 Health.

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

Abstract

Purpose: To develop a child self-report questionnaire measuring bilingual experience and self-perceptions of Spanish and English proficiency and establish preliminary evidence of validity and reliability for the questionnaire. **Method:** Participants included 113 Spanish-English bilingual children with and without developmental language disorders ranging in age from 4 to 8 years. All children completed the questionnaire in Spanish and participated in behavioral assessment of their language skills in both Spanish and English.

Results: Using confirmatory factor analysis, a model with three correlated factors (Self-Perception of Proficiency in Spanish, Self-Perception of Proficiency in English, and Bilingual Experience) emerged with the best global fit, reasonableness, consistency with theory, and model parsimony, suggesting that the questionnaire has good internal reliability. The scaled results of the questionnaires significantly correlated with behavioral measures of both Spanish and English, supporting the convergent validity of the measure.

Conclusion: The Houston Questionnaire is an assessment tool for the assessment of bilingual experience and self-perception of proficiency in Spanish and English bilingual children between the ages of 4 and 8 years. The results provide foundational evidence supporting the reliability and convergent validity of this tool.

Keywords: Bilingualism, assessment, self-report, child language

For the journal-formatted final version of this manuscript, see:

https://doi.org/10.1044/2022_JSLHR-21-00675

56 **Can Bilingual Children Self-Report their Bilingual Experience and Proficiency? The**
57 **Houston Questionnaire**

58 Bilingual children represent a heterogeneous group of children who vary in their bilingual
59 experiences and proficiency in each language (e.g., Bedore et al., 2010; Kapantzoglou et al.,
60 2015). This variation poses a significant challenge for the identification of language disorders in
61 bilingual children because speech-language pathologists must differentiate typical variations in
62 bilingual experience (e.g., children with less exposure to a language resulting in lower
63 proficiency in that language) from language ability limitations (e.g., language learning
64 difficulties; Arias & Friberg, 2017; Bedore & Peña, 2008). Therefore, it is critical to gather
65 information about the child's experiences in both languages during the bilingual assessment
66 process to better understand the potential impact of exposure and use on bilingual language skills
67 (Castilla-Earls et al., 2020; Kohnert, 2010).

68 Parents and teachers often serve as sources of information regarding the child's bilingual
69 experiences (e.g., Restrepo, 1998; Rojas et al., 2016). However, parents might be better at
70 estimating their child's abilities and experiences in the home language in comparison to the
71 school language. Parents in immigrant families may not speak the school language (National
72 Kids Count, 2020). Further, most parents do not have the opportunity to observe the child at
73 school, making it difficult to rate their school language use appropriately (Bedore et al., 2011).
74 Similarly, teachers might be limited in their ability to estimate children's language exposure and
75 use outside the school environment (Vagh et al., 2009). From this perspective, children
76 themselves might be better observers and reporters of their bilingual experience and knowledge
77 of each language than either parents or teachers. We developed The Houston Questionnaire
78 (Houston-Q) to gather information about bilingual experience and proficiency in Spanish and
79 English from the child's perspective.

80

81 **Self-Reporting of Bilingual Skills in Bilingual Children**

82 To develop a self-report measure of bilingual experience and proficiency, it is crucial to
83 first consider whether children have enough language awareness to express differences between
84 Spanish and English proficiency and experiences in each language. Language awareness is a
85 metalinguistic skill that requires the ability to reflect on one's own language (Svalberg, 2007).
86 Specifically for bilingual children, language awareness includes the ability to reflect on both of
87 their languages (Adesope et al., 2010). Language awareness in bilingual children develops as
88 early as age two. For example, two-year-old bilingual children can name their languages and
89 identify what language is being used by themselves and others (De Houwer, 2017).

90 Researchers examining language awareness in bilingual children have used various data
91 collection tools, including drawing and coloring language activities (e.g., color a child silhouette
92 following the languages spoken; Martin, 2012; Melo-Pfeifer, 2015; Rojo & Echols, 2017),
93 interviews (open questions about their bilingual experience that allow elaboration in responses;
94 Pérez-Leroux et al., 2011), and language questionnaires (Babino & Stewart, 2016; Rojo &
95 Echols, 2017). For example, Babino and Stewart (2016) used a 4-point Likert scale and multiple-
96 choice questions to examine cultural identity, language attitudes, and language use in and outside
97 the school. Language questionnaires emerged as an appropriate instrument, having been used
98 with bilingual children as young as four years of age (Rojo & Echols, 2017). In addition,
99 questionnaires allow for a variety of question types to elicit theoretical and practice-driven
100 information about language use, including yes/no questions (e.g., Do you use Spanish with your
101 teacher?), short open questions (e.g., Tell me a family member who lives in your house. What
102 language do you speak with him/her?), and quantifiable questions (e.g., How many friends do

103 you have who speak Spanish?). Questionnaires also can include visual aids to support more
104 reliable responses to quantitative questions. Therefore, a questionnaire appeared to be an
105 appropriate measurement tool for children to self-report their bilingual experience and
106 proficiency in each language.

107 Importantly for this study, the accuracy of children’s judgments of their bilingual
108 experience and proficiency has been largely unexplored. Previous studies investigating bilingual
109 children’s language awareness have primarily provided descriptive information about children’s
110 responses to the questionnaires (e.g., Babino & Stewart, 2016; Rojo & Echols, 2017). For a child
111 self-report questionnaire of bilingual experience and proficiency to be practically useful, it is
112 crucial to examine the descriptive information elicited by the tool and if children can respond in
113 a reliable and valid manner to the questionnaire. In this study, we aim to examine evidence of the
114 internal reliability and convergent validity of children’s responses to the Houston-Q, a self-report
115 questionnaire designed to quantify children’s bilingual experience and proficiency in each
116 language.

117 **The Development of the Houston Questionnaire**

118 The Houston-Q was designed to gather information about children’s self-assessment of
119 their language proficiency in both Spanish and English, and the child’s perceptions of their
120 bilingual experience. Other validated self-report measures exist for children to self-report similar
121 constructs (e.g., health-related quality life, stress, and psychological dysfunction; Pagano et al.,
122 2000; Solans et al., 2008; Osika et al., 2007). In bilingual adults, self-report studies show that
123 self-report measures of proficiency can be valid measurement instruments (e.g., LEAP-Q;
124 Marian et al., 2007). However, in some instances, mismatches between the classification of the
125 adult’s self-report of bilingual profile (Spanish dominant, English dominant, or Balanced) and

126 the adult's bilingual profile calculated from behavioral language measures have been reported
127 (Gollan et al., 2021; Tomoschuk et al., 2019). In this study, we focus on Spanish-English
128 bilingual children because they represent the largest bilingual population in the U.S., yet they
129 continue to be disproportionally represented in special education programs (Artiles et al., 2002;
130 Samson & Lesaux, 2009). Better understanding children's self-reported bilingual experience and
131 proficiency in each language may complement clinical assessment practices by facilitating
132 identification of children's baseline language experiences and strongest language prior to direct
133 comprehensive language assessment. A reliable indication of the child's strongest language
134 would be clinically meaningful in potentially reducing the time needed to problem-solve during
135 the bilingual evaluation process, particularly in the context of screening for language disorders.
136 Considering the child's abilities in their self-perceived strongest language may contribute to
137 more accurate identification of language disorders.

138

139 ***Bilingual Experience***

140 It is generally understood that exposure to a language is a prerequisite for language
141 learning and proficiency (e.g., Bohman et al., 2010; Hoff & Core, 2013). That is, for children to
142 learn a language, they need to be exposed to it. However, there is no agreement in the literature
143 about the amount and quality of the input needed for language learning (for a detailed review of
144 the methodological considerations regarding language input, see Carroll, 2017). For bilingual
145 children, language experiences are partitioned between two languages, in contrast with
146 monolingual children whose language input is completely in one language (Bridges & Hoff,
147 2012; Peña et al., 2018).

148 The amount of exposure a bilingual child has in each language robustly predicts their rate
149 of growth and proficiency in each respective language (e.g., Hammer et al., 2012; Hoff et al.,
150 2018; Peña et al., 2018). However, it is important to note that, in the U.S., English growth
151 predominates even among children who have high exposure to Spanish since exposure to English
152 tends to be greater outside the home, and Spanish exposure is likely to be limited to the home
153 (Hoff, 2017). On the other hand, Spanish exposure is necessary, although not sufficient, for the
154 development and maintenance of Spanish language skills of bilingual speakers, perhaps due to
155 the lower social status of Spanish in the U.S. (Castilla-Earls et al., 2019; Duursma et al., 2007).
156 Therefore, it is important to estimate how input is partitioned between languages to estimate
157 current exposure and potential future growth in each language.

158 An important part of the bilingual experience for children is the language(s) spoken at
159 home and its impact on language growth (De Houwer, 2004). For example, when both parents
160 speak Spanish at home, children tend to have higher vocabulary in Spanish than in English, but
161 when both parents speak English at home, English vocabulary tends to be higher than Spanish
162 vocabulary (Place & Hoff, 2011). Siblings also play a role in the bilingual experience at home.
163 For instance, homes with older school-age siblings tend to use more English than homes without
164 an older sibling (Bridges & Hoff, 2012; Obied, 2009). Interestingly, when bilingual college
165 students reflect on their experiences learning Spanish and English, they often attribute their
166 parents and grandparent's encouragement to use Spanish as an important contributor to their
167 current Spanish skills, while the use of English with siblings was considered a contributor to
168 their English skills (Castilla-Earls & Fulcher-Rood, unpublished).

169 The language(s) used at school also predicts language growth for children. For many
170 Spanish-English speaking children in the U.S., the start of formal education instigates a

171 significant shift in language proficiency from Spanish, the language spoken at home, to English,
172 the language spoken in most schools (Lutz, 2008). Children who attend bilingual education
173 schools tend to maintain Spanish language skills better than children who attend schools with
174 English-only instruction (Castilla-Earls et al., 2019; Farver et al., 2009). However, by 5th grade,
175 native Spanish-speaking children in bilingual education programs report that they prefer to use
176 English for both social and academic purposes (Babino & Stewart, 2016).

177 *Language Ability and Language Proficiency*

178 During the development of the Houston-Q, we aimed to capture the child's self-
179 assessment of language proficiency rather than language ability. In this study, language ability
180 refers to the child's general language learning capacity that interacts with language input (Peña et
181 al., 2018). Children with low language ability not explained by associated neurological disorders
182 are identified as children with developmental language disorders (DLD; Leonard, 2014; Bishop
183 et al., 2016). These children have low language ability even when input is present (Kan &
184 Windsor, 2010; Peña et al., 2014). Language ability is traditionally measured with standardized
185 language tests or spontaneous language measures (e.g., Peña et al., 2018; Restrepo, 1998). In
186 bilingual children, language ability is determined using the child's strongest language to
187 differentiate children whose language performance on a test or assessment task represent a lack
188 of input in a language (i.e., second language learners) from children who show low performance
189 in both languages (i.e., children with language disorders) (Kohnert, 2010; Peña et al. 2018).

190 Language proficiency refers to the specific knowledge of a language that is mediated by
191 the child's language ability. Regardless of whether a bilingual child has typical language ability
192 or low language ability, they will vary in their knowledge of each language. For example, a child
193 with low language ability may have more knowledge of Spanish than English, more knowledge

194 of English than Spanish, or have about the same level of knowledge of both languages. In the
195 same way, a child with typical language skills can vary in their bilingual profiles. However, how
196 much knowledge children with low language ability have in each language will differ from the
197 knowledge children with typical language ability have in their languages. That is, children with
198 low language ability that have about the same level of knowledge in both languages would score
199 lower on behavioral language assessments in comparison to children with typical language
200 ability who also have similar levels of knowledge in both languages. Therefore, there are at least
201 two levels of comparison¹. At one level, there is a between-child comparison of how much
202 language a child can learn provided input compared to their peers (language ability). At a second
203 level, there is a within-child comparison of how much knowledge a child has in a given language
204 compared to their other language (language proficiency). For the development of the Houston-Q,
205 we focused on this second level of comparison. We suggest that bilingual children can self-report
206 their proficiency in each language because they are aware of their two languages and can use
207 their awareness to respond to questions that yield to a proficiency or experience measure.
208 However, we do not suggest or expect that bilingual children would be able to self-report their
209 language ability (i.e., if they have a language disorder or typical language skills) because this is a
210 higher-level metalinguistic skill that requires a comparison between children.

211

212 **Measurement Reliability and Validity**

¹ There might be other levels of comparison, which are not the focus of this investigation. For example, a within child comparison of type of language skills, such as morphology and semantics (Bedore et al., 2012).

213 A core component of scale development is the evaluation of the scale's psychometric
214 properties, such as reliability and validity. This evaluation is an inherently ongoing process that
215 requires iterative examination of characteristics of the scale and how it functions for different
216 individuals in different contexts (see Boateng et al., 2018). In the present work, we focus on the
217 initial steps of psychometric evaluation, including examination of the developed measure's
218 dimensionality, the overall scale and subscale internal consistency reliability, and preliminary
219 convergent validity. These foundational properties directly influence the scoring structure and
220 interpretation of individual responses to a measure (American Educational Research Association,
221 American Psychological Association, & National Council on Measurement in Education, 2014)
222 and correspondingly represent a first step in establishing the practical utility of the Houston-Q.

223 Dimensionality assessment encompasses the identification of any potential subscales or
224 subtests within the overall measure. It is essential to establish the dimensionality of a measure
225 prior to evaluating its reliability because each unique dimension must be scored separately.
226 Scoring multiple dimensions together can lead to inaccurate estimates of item characteristics and,
227 ultimately, individual performance (de Ayala, 2013; DeMars, 2012; McNeish & Wolf, 2020).
228 Once subscales are identified, these can then be evaluated for evidence of internal consistency
229 reliability, which is the consistency within the test items included within each subscale. For a
230 subscale score to be meaningful, each test item included in that subscale should function in a
231 relatively similar manner. Internal consistency reliability is generally evaluated by examining
232 Cronbach's alpha or Coefficient omega in the case where some items contribute more strongly to
233 the total subscale score than others (i.e., violations of the assumption of tau equivalence;
234 McNeish, 2018).

235 Upon establishment of scale dimensionality and internal consistency reliability, evidence
236 of validity may be examined. Although there are many forms of validity, we focus on the
237 assessment of concurrent criterion validity, specifically convergent validity. Concurrent criterion
238 validity refers to how closely the scale and/or subscales are associated with scores obtained from
239 external measures administered to the same participants at approximately the same time
240 (American Educational Research Association, American Psychological Association, & National
241 Council on Measurement in Education, 2014). Convergent validity may be evaluated empirically
242 through the examination of correlations between participants' scores on the developed scale and
243 their scores on other measures that are theoretically considered to be related. Examination of
244 these psychometric properties goes beyond a descriptive approach to the responses provided by
245 children (e.g., Martin, 2012; Melo-Pfeifer, 2015; Rojo & Echols, 2017) and instead targets the
246 quality of the measurement.

247

248 **This study**

249 Previous research suggested that bilingual children as young as four years of age might
250 have enough language awareness to self-report their bilingual experiences and proficiency in
251 each language (Rojo & Echols, 2017). However, information about children's bilingual
252 experience and proficiency is currently collected primarily through parents and teachers. In this
253 study, we explore the possibility that children can provide a valid and reliable estimation of their
254 bilingual experiences and proficiency using a questionnaire administered verbally in Spanish by
255 an adult. We developed the Houston-Q as a tool to estimate variations in bilingual experience
256 and proficiency during the bilingual assessment process. Our research questions are: (a) what is
257 the dimensionality of the Houston-Q? (b) Is the Houston-Q a reliable tool for the self-report of

258 bilingual experience and proficiency in Spanish and English in bilingual children? and (c) Is
259 there evidence of convergent validity between behavioral measures of language skills in Spanish
260 and English and the Houston-Q?

261

262

Method

263 Validation Participant Sample

264 The Institutional Review Board at the University of Houston approved this study. Parents
265 provided written informed consent, and children provided verbal assent to participate in the
266 sessions. Participants were recruited from school districts and speech-language clinics in the
267 Greater Houston area as part of a broader longitudinal study of bilingual language development.
268 To be eligible for the study, children spoke and understood both Spanish and English, passed an
269 otoacoustic emission hearing screening, and obtained a score greater than 70 on the Matrices
270 subtest of the *Kaufman Brief Intelligence Test—Second Edition* (KBIT-2; Kaufman & Kaufman,
271 2004) as a measure of non-verbal IQ².

272 The validation sample for the current study included 113 Spanish-English bilingual
273 children ranging in age from 3 years, 11 months to 8 years, 2 months ($M = 71.05$, $SD = 12.46$ in
274 months). The sample was 43% girls ($n = 49$). Approximately 54% of the children came from
275 families where the mother had not attended college, and 70% of the children qualified for free or
276 reduced-price lunch as reported via parental questionnaire. Parents also reported that their
277 families spoke either Spanish only (49%) or both Spanish and English at home (39%). The
278 remaining 12% of the parents reported their children spoke English only at home and Spanish at

² There was one instance of a child with a score below 70 on the KBIT-2 but with all scores on language assessments within normal limits. It appears that the KBIT-2 score was not indicative of the child's actual abilities. For this reason, we ran all analysis twice: a) excluding this child, and b) including this child. We found no differences in the results of this study. Therefore, we included this child in the reported sample.

279 school. Regarding the language of instruction at school, 90% of the children in our sample
280 attended bilingual Spanish-English or Spanish language immersion education programs. Further
281 information about the children's language skills will be presented as descriptive information in
282 the results section.

283

284 **Measures**

285 *The Houston Questionnaire*

286 The Houston-Q was developed to provide a self-assessment measure for children
287 regarding their language proficiency and experiences in each language. The questionnaire was
288 constructed to be completed in approximately 10 minutes with children as young as four years
289 old. For this reason, we designed questions with simple wording and vocabulary and used visual
290 support when needed. In addition, all questions were designed to be verbally presented in
291 Spanish by an examiner who recorded the child's responses. Questions included yes/no
292 questions, short open questions, and questions with Likert-scale options to obtain quantifiable
293 information. Some questions required a combination of yes/no responses followed by a 5-point
294 Likert-scale question (e.g., 1- *a little* to 5- *a lot*; 1- *few* to 5- *many*). To support children, we used
295 pictures with different amounts of candy to indicate a little or few (one piece of candy) to a lot or
296 many (five pieces of candy). Other questions asked about home and school activities and the
297 language in which they occurred (Spanish, English, both, or not performed at all).

298 The questionnaire was designed to target three main areas of children's language: self-
299 assessment of Spanish proficiency, self-assessment of English proficiency, and bilingual
300 experience in Spanish and English. It consists of 25 questions in total. The section for self-
301 assessment of proficiency in the languages in the Houston-Q includes questions regarding how
302 good children are at speaking a language, how easy they perceive the language to be, and how

303 many friends they have who speak the languages. These questions are a combination of a yes/no
304 question (e.g., Are you good at speaking Spanish? Do you think Spanish is easy? Do you have
305 friends who speak only Spanish?) followed by a 5-point Likert-scale question (e.g., If you are
306 good, how good? If it's easy, how easy? If you have friends who speak that language, how
307 many?). On the 5-point Likert scale follow-up questions, lower values indicated lower
308 proficiency, and higher values indicated higher proficiency. Other items in the proficiency
309 section of the Houston-Q included questions regarding how much Spanish and English children
310 heard during the day, which were also 5-point Likert scale questions with lower values indicating
311 lower frequency and higher values indicating higher frequency. To estimate bilingual experience,
312 questions listed a variety of activities (e.g., read books, watch TV, play at the park, etc.) and
313 children were provided four options regarding the language they used during these activities
314 (e.g., I do this in Spanish, I do this in English, I do this in both Spanish and English, and I don't
315 do this). A final set of questions prompted children to name three people from their family and
316 identify what language they used with each of them. Children were provided with three options
317 to respond: Spanish, English, or both Spanish and English.

318

319 *Behavioral language measures*

320 **Receptive Vocabulary.** We used the standard scores from the *Peabody Picture*
321 *Vocabulary Test—Fourth Edition* (PPVT-4; Dunn & Dunn, 2007) and the *Test de Vocabulario*
322 *en Imágenes Peabody* (TVIP; Dunn et al., 1986) as measures of receptive vocabulary in English
323 and Spanish, respectively. The PPVT-4 is a standardized measure of receptive vocabulary for use
324 with individuals ages 2-90 years old. This assessment has been normed with English
325 monolinguals from across the United States and has been frequently used in research studies with

326 children as a measure of vocabulary. The TVIP is a parallel measure to the PPVT and assesses
327 receptive vocabulary in Spanish in individuals ages 2-18. The TVIP has been normed with
328 Spanish monolingual speakers in Mexico and Puerto Rico. In both assessments, children are
329 presented with stimulus pages consisting of four pictures. The examiner provides a vocabulary
330 word to the child, and the child responds by either pointing or stating the number for the picture
331 they believe best represents the word. It is important to note that both of these tools were normed
332 with monolingual children and therefore are not ideal for measuring receptive vocabulary
333 abilities in bilingual children (Wood et al., 2018).

334 **Morphosyntax.** We used the morphosyntax subtests of the *Bilingual English-Spanish*
335 *Assessment* (BESA; Peña et al., 2018) and the *Bilingual English-Spanish Assessment—Middle*
336 *Extension* (BESA-ME; Peña et al., 2008, 2016). The BESA is a standardized test designed to
337 evaluate the language abilities of Spanish-English bilingual children ages 4;0-6;11 (years;
338 months) in the U.S. The BESA-ME is an experimental measure, similar to the BESA, to assess
339 language skills of Spanish-English bilingual children ages 7-9;11 (years; months). The BESA
340 (and BESA-ME) was used in this study to estimate language ability because it is currently the
341 gold standard normed-reference measure for identification of Spanish-English bilingual children
342 with developmental language disorders in the United States. The morphosyntax subtest of both
343 tests consists of a cloze item section and a sentence repetition section targeting complex
344 grammatical structures in each language. Standard scores ($M = 100$, $SD = 15$) are calculated for
345 each language. The BESA and BESA-ME morphosyntax subtests can be administered as stand-
346 alone subtests with good diagnostic accuracy to identify bilingual children with developmental
347 language disorders (Peña et al., 2008, 2016, 2018). In order to combine BESA and BESA-ME
348 morphosyntax, we used standard scores. The BESA morphosyntax subtests standard scores range

349 from 52-145. However, the BESA-ME experimental version standard scores did not have a
350 specific range at this time. For purposes of the analyses in this study, we mirrored the range on
351 the BESA-ME to the one used for the BESA so that the lowest possible score on the BESA-ME
352 was also 52³. We used the best language score as a measure of language ability, as suggested in
353 the BESA and BESA-ME testing manuals, following current best practices for the assessment of
354 bilingual children (Kohnert, 2010; Peña et al., 2018).

355 **Sentence Repetition.** We also used the scaled scores of the Recalling Sentences subtest
356 (Recordando Oraciones in the Spanish version) in the latest versions of the *Clinical Evaluation*
357 *of Language Fundamentals* in English and Spanish (CELF-5 for English, Wiig et al., 2013; and
358 CELF-4 for Spanish, Semel et al., 2006). In these subtests, children are asked to repeat the
359 sentence after the evaluator. The subtest is designed to evaluate the child’s knowledge of the
360 language structure and vocabulary in addition to cognitive processing skills such as verbal
361 working memory (Pratt et al., 2020). Because this task assesses the knowledge of the language
362 (i.e., to be able to repeat a sentence, one needs to have the language structure and vocabulary in
363 that language), sentence repetition tasks might be considered biased for the assessment of
364 language ability in bilingual children if only one language is used (Armon-Lotem & Meir, 2016).
365 Sentence repetition tasks have been found to have high sensitivity and specificity for the
366 diagnosis of developmental language disorder (Archibald & Joanisse, 2009; Rujas et al., 2021).

367 **Procedures**

368 Parents provided consent for their children’s participation in the study and completed a
369 questionnaire about demographics and the use of Spanish and English. Children provided assent
370 to participate. Children completed the behavioral language tasks and the Houston-Q as part of a

³ Sensitivity analyses were performed to assess the potential impact of the truncated scores on correlational results. No substantial differences were noted, so only the results from the truncated scores are reported.

371 larger battery of assessments. The Spanish language tasks were part of the Spanish language
372 skills session, and the English language tasks were part of the English language skills session.
373 Each of these sessions was approximately 50 minutes long. Task order in each session varied
374 across participants. All the tasks were administered in person and scored by a trained research
375 assistant who was a native speaker of the target language.

376 The Houston-Q was administered in Spanish as part of the Spanish language skills
377 session. Children were first shown pictures with different amounts of candy to indicate a little
378 (one piece of candy) to a lot (five pieces of candy). The examiner said in Spanish, “I know you
379 speak both Spanish and English; I am going to ask you some questions about Spanish and
380 English. For some questions, you can answer a little, like one piece of candy; for others, you can
381 answer a lot, like five pieces of candy. For some questions, you may want to answer something
382 in between, like two, three, or four pieces of candy.” The examiner gauged the child’s
383 understanding of the task by asking questions to ensure that the child understood what was
384 expected (e.g., Do you have any questions? Do you understand what we are doing?). Once the
385 examiner felt that the child understood the task, they would start asking the questions in the
386 Houston Questionnaire. The examiner monitored whether the child answered each question in a
387 manner aligned with the intended content of the question to ensure understanding of the task.
388 Repetition of the instructions was allowed. The examiner wrote down all answers from the child
389 in the questionnaire response form. Although the questionnaire was administered in Spanish
390 only, responses were allowed in Spanish or English. All children in this study were able to
391 complete this task using this procedure. There were no reports of no compliance or difficulties
392 understanding the task.

393

394 **Analytic Approach**

395 Children's responses were first examined for frequencies of each response (see
396 Supplementary Figures 1-3). Item responses were evaluated for evidence of restriction of range
397 (i.e., floor or ceiling effects), which would limit information extractable from any given item,
398 based on a criterion of 95% for any specific response. No items met this criterion.
399 Correspondingly, all items were included in subsequent analyses.

400 ***Dimensionality and Reliability***

401 We used confirmatory item-level factor analysis to assess the dimensionality, or
402 underlying factor structure, of the scale. An inherent strength of this analytic approach is that it
403 allows for the evaluation of the characteristics of individual questionnaire items by partitioning
404 out different sources of variability in children's responses. Item-based confirmatory factor
405 analysis yields separate estimates for individual item characteristics (e.g., difficulty,
406 discrimination) and individual participant characteristics (e.g., self-perception of Spanish
407 proficiency). This analysis is useful for supporting the development of a generalizable scale.
408 However, the robustness of the specific item parameters is limited by the representativeness of
409 the participant sample compared to the local population.

410 We based all model testing on a priori hypotheses of possible constructs underlying the
411 items. The most complex model assessed included six possible underlying factors (see Figure 1,
412 Model A), and the most parsimonious included three underlying factors (Figure 1, Model B), in
413 alignment with the construction of the scale. All factors were correlated, consistent with the
414 theoretical framing that general language learning abilities contribute to the development of
415 proficiency in both languages. Models were estimated using unweighted least squares means and
416 variance (ULSMV) in Mplus Version 8.4 (Muthén & Muthén, 2019). Item intercepts, factor

417 loadings, and residual variances were freely estimated, with latent factor means fixed at 0 and
418 latent factor variances fixed at 1 for model identification.

419 Model fit was assessed through (a) evaluation of parameter estimates and residuals, with
420 models examined for evidence of misfit through indicators such as negative residual variances
421 and unexpectedly large or small estimates; (b) consideration of global fit indices, including the
422 chi-square test of model fit, root mean square error of approximation (RMSEA), comparative fit
423 index (CFI), Tucker Lewis index (TLI), and standardized root mean square residual (SRMR)
424 following guidance summarized by Lomax (2013); and (c) chi-square difference testing of nested
425 models using the DIFFTEST option for ULSMV in Mplus (Muthén & Muthén, 2019). More
426 parsimonious models were favored when no significant difference in global fit was observed.

427 There were two items that, from a theoretical perspective, could contribute to more than
428 one underlying construct. These items were #10 “¿Tienes amigos que hablen inglés y español? /
429 Do you have friends who speak English and Spanish?” and the follow-up question #11
430 “¿Cuántos? / How many?” We hypothesized that these two items might reflect Spanish exposure
431 and English exposure, or they might only reflect to Spanish exposure (because English is the
432 majority language in the U.S.). To assess this, we compared models including these items cross-
433 loaded onto both factors to models with the items only loaded onto the Spanish exposure factor.

434 Upon identification of the underlying structure with the best balance of model fit,
435 parsimony, and alignment with theoretical construction, we computed reliability indices for each
436 subscale identified. Coefficient omega hierarchical was used to accommodate potential
437 violations of tau equivalence (McNeish, 2017).

438 ***Practical Scoring Approaches***

439 We considered several scoring approaches for practical use of the scale, drawing on
440 related discussion from DiStefano et al. (2009) and Logan et al. (2019). Ease of administration
441 and interpretation is essential to the practical, day-to-day useability of assessments.
442 Consequently, we examined a restriction on the factor loadings, which required each item to
443 contribute equally to its corresponding subscale. This analysis is similar to comparing a 2-
444 parameter item response theory (2-PL IRT) model to a 1-PL IRT model. We compared global fit
445 for the restricted model to a model without restriction. We also obtained metrics of parameter
446 bias to determine the practical difference between equal weighting of items compared to
447 differential item weighting within each subscale. Based on the results, we constructed a
448 preliminary useable system for scoring the measure.

449 ***Convergent Validity***

450 After identifying the underlying structure with the best fit to the data, we examined
451 indicators of convergent validity for the Houston-Q. To do this, we used the developed measure
452 to compute scores for each scale construct for each child. We then examined correlations among
453 the obtained scale scores and concurrent measures of Spanish and English language. The
454 concurrent measures of language in Spanish were BESA/BESA-ME Morphosyntax, CELF-4
455 Recordando Oraciones, and TVIP. In English, the three language measures were BESA/BESA-
456 ME Morphosyntax, CELF-5 Sentence Recall, and PPVT-4. We expected the subscales of self-
457 reported Spanish proficiency to be positively associated with the Spanish language measures and
458 the subscales of self-reported English proficiency to be positively associated with the English
459 language measures. Similarly, we hypothesized that the subscales of bilingual experience would
460 correlate with the Spanish and English measures, such that greater Spanish experience would

461 correspond with higher Spanish language scores and greater English experience would
462 correspond with higher English language scores. Finally, we examined correlations between
463 children's subscale scores on the Houston-Q and age.

464

465

Results

466

Descriptive information

467

468

469

470

471

472

473

474

475

476

477

478

479

Children in the sample varied widely in terms of their language proficiency profiles. To illustrate this variability, we descriptively examined participants' standard scores on the language measures used in this study separately for each language. These included standard scores in Spanish and English for the BESA/BESA-ME, sentence repetition subtest of the CELF-4 in Spanish and CELF-5 in English, and receptive vocabulary using the PPVT and TVIP. For 46% of the children in this sample, the difference between their scores in Spanish and English for the BESA/BESA-ME were within 10 standard points of each other, suggesting that about half of the children had relatively balanced morphosyntactic skills in both languages. For the remaining children, 31% had stronger English morphosyntactic skills (more than a 10-point difference in standard scores), while 23% had stronger Spanish skills. For vocabulary, 41% of the children had scores in Spanish and English within ten standard points of each other. In comparison, 24% of the children had stronger receptive vocabulary in English, and 35% had stronger Spanish receptive vocabulary.

480

481

482

483

Children in this sample also varied in language ability. The average score in the best language for the BESA/BESA-ME was 92.19 ($SD = 17.06$), for PPVT was 82.63 ($SD=25.00$), and for TVIP was 84.53 ($SD=25.50$). Forty-two percent of the children were receiving speech-language services in their schools. These aspects of language ability, proficiency, and use

484 indicate that our participants represent a heterogeneous group of bilingual children. Detailed
485 information for the children in our sample is included in Table 1.

486 **Sample Characteristics**

487 Response frequencies for each item in the questionnaire are depicted in the
488 Supplementary Material (Figures S1-S3). Generally, children in the present sample rated
489 themselves as speaking both Spanish and English well (Spanish, $n = 101$, and English, $n = 103$,
490 out of 113). However, the degree of how well children rated themselves as speaking each
491 language varied. Children were slightly more likely to report Spanish as being easy ($n = 97$ out
492 of 112) than English being easy ($n = 87$ out of 112), with more variability present in the reported
493 degrees of English easiness compared to Spanish (Figure S1).

494 On items focused on bilingual experience, children were asked about different activities
495 and whether these were done using both Spanish and English, only Spanish, or only English.
496 There was also an option to indicate that they did not do the activity. Of these options, children
497 most often reported using both languages during the activities. With their classroom teacher,
498 60% of children indicated that they used both Spanish and English, 21% used only Spanish, and
499 19% used only English. With respect to reading books, 59% of children reported reading in both
500 languages, 28% read only in Spanish, and 13% read only in English. Similarly, 66% reported
501 learning to write in both languages, 22% reported learning to write only in Spanish, and 13%
502 reported learning to write only in English. When watching TV, 54% of children watched in both
503 languages, 16% watched only in Spanish, and 30% watched in only English. At the park, 41%
504 played using both Spanish and English, 29% used only Spanish, and 31% used only English. In
505 family reunions, 41% of children used both languages, 34% used only Spanish, and 25% used
506 only English. These findings are provided in Figure S2.

507 Overall, children reported having both family members and friends who spoke Spanish,
508 English, and a combination of Spanish and English. When asked how much Spanish and English
509 they heard each day, 27% of the children reported hearing a lot of both Spanish and English.
510 Sixty-three children out of 112 reported hearing a lot of Spanish per day. Finally, 54 children out
511 of 112 reported hearing a lot of English per day (Figure S3).

512 No patterns were observed in missing data. Children elected not to respond to questions
513 randomly, with 32 instances of missing responses. Given that 3,051 total responses were possible
514 (27 items and 113 total participants), and no patterns were observed, data were considered
515 missing at random. We also examined patterns in children's responses for evidence of
516 contradictory patterns or illogical response combinations. The questionnaire items were written
517 to allow for all possible combinations of responses, but one noteworthy pattern occurred among
518 12 participants. Six children indicated that they were not good at speaking English but thought
519 English was easy. Another six children stated that they were not good at speaking Spanish but
520 thought Spanish was easy. Although this combination of perceptions seems unlikely, individuals
521 can have the belief that learning a language is easy, even though they do not consider themselves
522 to be good at speaking that language. Consequently, we did not interpret these response
523 combinations as problematic.

524 **Dimensionality and Reliability**

525 Confirmatory factor analyses indicated that a model with three correlated factors yielded
526 the best balance of global fit, reasonableness, consistency with theory, and model parsimony (see
527 Figure 2). The model included a single factor underlying the items designed to measure
528 children's self-perceptions of their proficiency in Spanish (i.e., "Self-Perception of Spanish"), a
529 single factor underlying items designed to measure children's self-perceptions of their

530 proficiency in English (i.e., “Self-Perception of English”), and a single factor underlying self-
531 reported bilingual experience (i.e., “Bilingual Experience”). This model, specified with item
532 loadings and thresholds freely estimated, provided a good fit to the data: $\chi^2(296) = 325.90$ and p
533 $= .1118$, RMSEA = 0.030 (90% CI [0.001, 0.048]), CFI = 0.936, TLI = 0.930, SRMR = 0.114.
534 Coefficient omega hierarchical was computed to be .910 for Self-Perception of Spanish, .753 for
535 Self-Perception of English, and .893 for Bilingual Experience indicating that the three factors
536 showed good internal consistency reliability.

537 The two items that were hypothesized to contribute to more than one underlying factor
538 (#10 “¿Tienes amigos que hablen inglés y español? / Do you have friends who speak English and
539 Spanish?” and follow-up question #11 “¿Cuántos? / How many?”) were examined as indicators
540 of Self-Perception of Spanish and of Self-Perception of English. Item loadings and model
541 comparisons suggested that item #10 did not fit well on either factor, whereas #11 contributed
542 reasonably to children’s Self-Perception of Spanish. Chi-square testing of model B (see Figure 1)
543 with item #10 freely loaded onto Self-Perception of Spanish compared to being fixed at zero
544 resulted in no significant difference in fit: $\Delta\chi^2(1) = 0.10$, $p = .751$. Item #10 was removed from
545 subsequent modeling, and #11 was loaded onto only the Self-Perception of Spanish proficiency
546 factor. Global model fit statistics and chi-square comparisons of nested models are provided in
547 Table 2. Standardized item loadings and thresholds are provided by item in Table 3.

548 **Scoring System for the Houston-Q**

549 When item loadings were restricted to be equivalent (analogous to a 1-PL IRT model),
550 global model fit comparisons revealed a significantly worse fit to the data compared to the model
551 with freely estimated loadings $\Delta\chi^2(23) = 59.34$, $p < .001$. Additionally, this restriction resulted
552 in a total parameters bias of 35% across the subscales, with the least bias observed for the

553 Bilingual Experience factor (28%) compared to the Self-Perception of English (41%) or Self-
554 Perception of Spanish (40%) factors. Consequently, the free estimation of item loadings was
555 retained for the preliminary scoring system of the measure, which was constructed based on the
556 standardized weighted contributions of each item (see Houston-Q Español, Houston-Q English,
557 and Houston-Q Research spreadsheets). Given the random missing data patterns observed in the
558 data used for the present study, the scoring system is designed to allow for the computation of
559 scores with missing individual item responses.

560 The measure was scaled from 0-10 for the Self-Perception Scores of Spanish and English
561 proficiency, where 0 = no proficiency and 10 = full proficiency. For Bilingual Experience, we
562 scaled responses from 0-20, with 0 = all experiences in Spanish, 10 = equal experiences in
563 Spanish and English, and 20 = all experiences in English. We elected to scale the values
564 differently to reflect the differences in the underlying constructs.

565 **Convergent Validity**

566 Within the present participant sample, children scored an average of 7.73 ($SD = 2.15$) for
567 Self-Perception of Spanish proficiency, suggesting relatively high levels of self-perceived
568 proficiency in Spanish. Self-Perception of English was similarly high, with an average of 7.69
569 ($SD = 2.09$). The children's self-perception scores for proficiency in each language were
570 significantly and positively associated with the behavioral measures of language with small to
571 moderate correlations. The self-perception scores were negatively associated across languages (r
572 = $-.24$, 95% CI [$-.40$, $-.05$], $p = .013$), indicating that children who reported high proficiency in
573 Spanish tended to report lower proficiency in English and vice versa. Self-Perception of Spanish
574 correlated with the Spanish measures CELF-4 Recordando Oraciones, TVIP, and BESA
575 Morphosyntax at $r = .36$ (95% CI [$.19$, $.51$], $p < .001$), $r = .23$ (95% CI [$.04$, $.40$], $p = .017$), and

576 $r = .42$ (95% CI [.25 .56], $p < .001$), respectively. Self-Perception of English similarly correlated
577 with the English measures CELF-5 Sentence Repetition, PPVT-4, and BESA Morphosyntax at r
578 $= .32$ (95% CI [.14, .48], $p < .001$), $r = .24$ (95% CI [.05, .40], $p = .013$), and $r = .23$ (95% CI
579 [.04, .40] $p = .017$), respectively.

580 On average, children indicated generally balanced bilingual experience, with slightly
581 higher experience in Spanish than English, evidenced by the average Bilingual Experience at
582 8.94 ($SD = 4.53$). Appropriately, increased experience in Spanish was associated with a higher
583 self-perception of Spanish proficiency: $r = -.61$ (95% CI [-.72, -.48] $p < .001$), and increased
584 experience in English was associated with a higher self-perception of English proficiency: $r =$
585 $.42$ (95% CI [.25, .56] $p < .001$). Age correlated weakly with self-perception of Spanish ($r = -.19$,
586 95% CI [-.36, -.01], $p = .050$), but not with the other two subscales. See Table 4 for full
587 correlations.

588

589

Discussion

590 This study aimed to examine the reliability and convergent validity of the Houston
591 Questionnaire in a sample of young bilingual children. In this study, we included children with
592 varying levels of bilingual proficiency and language ability to capture variability in bilingual
593 experiences and proficiency. Our results provide initial evidence supporting the internal
594 consistency reliability and preliminary criterion validity of the Houston Questionnaire as a child
595 self-report assessment tool.

596 Dimensionality and Reliability

597 Our findings indicate that three correlated factors underlie children's responses to the
598 Houston-Q: Self-Perception of Spanish Proficiency, Self-Perception of English Proficiency, and

599 Bilingual Experience. These three factors were moderately correlated, which suggests that
600 participants' responses reflected distinct but related constructs. Each subscale had overall good
601 internal consistency reliability, which indicates that the questionnaire items were generally
602 cohesive within each factor (Revelle & Condon, 2019). These results suggest that Houston-Q can
603 elicit reliable responses from young bilingual children. In other words, the questions of the
604 Houston-Q elicit responses that are generally consistent in terms of bilingual experience and self-
605 ratings of Spanish and English proficiency. For example, a child is likely to respond that they are
606 good at speaking Spanish and that Spanish is easy. This appropriate internal consistency
607 reliability is crucial for a self-report measure since the questions must reliably measure the same
608 construct (Dunn et al., 2014; McNeish, 2017). Failing to do so would suggest that the measure is
609 not designed appropriately (e.g., not worded properly) or that different constructs are being
610 measured (e.g., constructs other than bilingual experience).

611 The questionnaire items aligned well with the hypothesized underlying factors. For
612 example, the questions that we expected to reflect Self-Perception of Spanish Proficiency were
613 reliably associated with one another. The same was found for Self-Perception of English
614 Proficiency and Bilingual Experience. There was no evidence of misfit in the final model, which
615 suggests that obtaining these three subscale scores from the Houston-Q is appropriate.

616 We hypothesized that two questionnaire items could contribute to Self-Perception of
617 Spanish proficiency and/or to Self-Perception of English proficiency. We directly tested the fit of
618 question #10, "Do you have friends who speak Spanish and English?" and follow-up question
619 #11, "How many?" as indicators of these underlying factors. The results indicated that question
620 #10 did not directly align with either self-perception of Spanish or self-perception of English, but
621 question #11 did align with self-perception of Spanish. We interpret these findings as primarily

622 reflective of the sampling context in Houston. In the present sample of participants, most
623 children reported having at least some friends who speak Spanish and English, which resulted in
624 relatively limited variability (i.e., restriction of range) for question #10. This limited variability
625 restricted the item's potential to contribute to any factor. Question #11, however, did result in
626 sufficient response variation to serve as an indicator of self-perception of Spanish. Because
627 English is the predominant language used in the U.S. and especially in schools in the U.S., it is
628 reasonable that children who report having more friends who speak both English and Spanish
629 would similarly have a greater self-perception of their Spanish proficiency.

630 For the present study, we purposefully included children with diverse ranges of exposure
631 and from a relatively broad age range to reflect the variability typically seen among bilingual
632 children in the U.S. However, these bilingual children are speakers of Spanish in a city where
633 Spanish is frequently heard and used by the broader community, and where opportunities for
634 formal education in Spanish exist. Therefore, these results provide initial evidence supporting the
635 utility of the Houston-Q across these characteristics. If there were substantial differences in the
636 validity or reliability of the measure between the subgroups, we would expect evidence of lack of
637 fit such as poor global model fit and spurious parameter estimates. Instead, we found that the
638 global fit of the model was good, especially given the relatively small sample size, and the
639 parameter estimates were generally stable. Although replication is certainly necessary to further
640 explore the validity, reliability, and overall functioning of the scale across subpopulations of
641 bilingual learners, the current findings provide preliminary evidence of the utility of the scale
642 across diverse Spanish-English speaking learners.

643 **Scoring System for the Houston-Q**

644 Using the results of the confirmatory factor analysis, we created scaled scores for Spanish
645 and English self-perception of proficiency and bilingual experience. A scale of 0 to 10 was used
646 to describe Self-Perception of Proficiency in Spanish and English and a scale from 0 to 20 to
647 describe Bilingual Experience. Importantly, we weighted the contribution of each item within
648 each scale to align with its unique factor loading, given that the items did not equally reflect the
649 underlying constructs of interest. We tested whether the items could be scaled to contribute
650 equally but found that this significantly worsened the reliability of the questionnaire. Forcing the
651 items to contribute equally resulted in substantial bias (i.e., 28 – 41%) in each subscale score. In
652 other words, weighting items equally resulted in substantially different subscale scores when
653 compared to varying the item weights. These results suggest that some of the questionnaire items
654 were more important indicators of children’s self-perception proficiency and bilingual
655 experience than others. For example, question #1, “Are you good at speaking Spanish?” was a
656 more robust and consistent indicator of self-perception of Spanish proficiency across children
657 than question #6, “Do you have friends who only speak Spanish?”. For question #1, the response
658 “yes” reliably reflected a higher overall self-perception of proficiency in Spanish. Children who
659 received a high score on self-perception of proficiency in Spanish generally responded “yes” to
660 question #1. On the other hand, for question #6, more friends who speak Spanish typically but
661 not always reflected higher self-perception of Spanish proficiency. There was a weaker
662 association between children’s total scores for self-perception of Spanish proficiency and their
663 responses to question #6. This variation in item contributions was observed for all three
664 subscales and is evident in the standardized item loadings. Our scoring system reflects this
665 variation by weighting each item differently.

666 The scoring system also allows children to receive scores on each of the subscales even if
667 they do not respond to individual questionnaire items. We incorporated this design feature
668 because the results of the present work revealed no patterns in children’s missing data,
669 suggesting that children randomly skipped questions throughout the questionnaire. Children did
670 not frequently skip items, and when they did, there was no apparent reason why they skipped.
671 We believe this may be attributable to normal lapses in attention. Consequently, it is reasonable
672 to obtain a subscale score even when children skip a few items across the questionnaire.

673

674 **Convergent Validity**

675 Children’s self-perception of Spanish proficiency correlated positively with the Spanish
676 language measures. Correlations with sentence repetition and productive morphology were
677 moderate, and correlations with receptive vocabulary were weak-to-moderate. Similarly, self-
678 perception of English proficiency positively correlated with the English language measures
679 overall. In English, the correlations with sentence repetition, receptive vocabulary, and
680 productive morphology were weak-to-moderate. Although replication with an independent,
681 larger sample is necessary to establish the magnitude of these associations more definitively, the
682 direction of the correlations is consistent. It is important to note that the receptive vocabulary
683 measures were normed on monolingual children and, therefore, are not appropriate estimation of
684 the vocabulary knowledge of the bilingual children in this study, which may have lowered the
685 magnitude of the correlations between the Houston-Q subscale scores and vocabulary,
686 particularly for Spanish.

687 As expected, children’s bilingual experience scores on the Houston-Q, which ranged
688 from 0 to 20, with 10 indicating fully balanced experience in Spanish and English, also

689 correlated with the external standardized measures. Bilingual experience correlated moderately
690 positively with self-perception of proficiency by language. Bilingual experience values between
691 0 and 10, which indicate more self-reported experience in Spanish, generally corresponded with
692 higher Spanish language scores. Further, bilingual experience values between 10 and 20, which
693 indicate more self-reported experience in English, generally corresponded with higher English
694 scores. These results suggest the bilingual experience metric functions as expected, with self-
695 reported exposure and use to each language aligning with norm-referenced scores in each
696 respective language.

697 We interpret these small to medium correlations and the direction of the associations to
698 be good indicators of the validity of the Houston-Q (Strauss & Smith, 2009). These correlations
699 indicate that proficiency measures using behavioral tasks and the children's perception of their
700 proficiency shared some properties, but they represent distinct constructs. This finding might be
701 explained by the fact that the behavioral tasks tap into specific language skills, like children's
702 ability to recall sentences, which might not necessarily be what children consider would qualify
703 them as good speakers of a language. Because we did not design the study a priori to compare
704 the strength of the correlations, we cannot make specific claims about what correlations are
705 stronger or weaker than others. However, the directionality of the correlations provides us with
706 necessary evidence informing the validity of Houston-Q. Recall that children's self-perception of
707 Spanish proficiency using the Houston-Q correlated positively with receptive vocabulary,
708 productive morphosyntax, and sentence repetition in Spanish while correlating negatively with
709 the same measures in English. Namely, children who rate themselves as good speakers of a
710 language tend to have higher scores from behavioral tasks in that language than children who
711 consider themselves not to be good at speaking that language. Further, children who rated

712 themselves as high in both languages tended to have high scores in both languages. These
713 findings suggest that children's responses to the Houston-Q rating are tapping into their
714 proficiency in each language.

715

716 **Sample-Specific Considerations**

717 It is important to consider that most children in this study rated themselves as speaking
718 both Spanish and English well, although the degree of their ratings varied. This consideration is
719 particularly important because about 40% of the children in this study had standard scores for
720 morphosyntax (BESA/BESA-ME) and receptive vocabulary (PPVT/TVIP) within 10 points of
721 each other, which suggests that their proficiency in each language was at similar levels. These
722 data need to be interpreted within the context. The data for this study was collected in Houston, a
723 city where 39.3% of the overall population speak Spanish at home (U.S. Census Bureau, 2019).
724 Bilingual education is available for children with limited English language ability because Texas
725 law mandates bilingual instruction for elementary schools with 20 or more children who need
726 English language support (Bilingual Education and Training Act). Spanish immersion is also
727 available in some schools in Houston, but it is not mandated by law. Notably, a significant
728 proportion of children in our sample attended bilingual programs and immersion programs. This
729 strong bilingual context might impact the child's ability to rate themselves in each language
730 since their everyday experiences include both Spanish and English, which might be different
731 from other contexts in the U.S. Therefore, future studies should be conducted in other bilingual
732 populations to examine the effect of the context on the reliability and validity of the Houston-Q.

733 The finding that self-perception of Spanish proficiency was associated with age and our
734 operationalized metric of language ability is worth noting. The shift into more English-focused

735 environments as bilingual children get older in the U.S. may explain the negative relationship
736 between age and participants' self-perception of Spanish proficiency. Recall that we included
737 children between 4 and 8 years of age in this study. At age 4, children tend to spend more time in
738 the home with their family, whereas by age 8, they are likely spending more time in the
739 community and with friends. Although bilingual education offers a protective effect on the
740 maintenance of Spanish language skills, it is not sufficient for some children (Castilla-Earls et
741 al., 2019). This interpretation is supported by the finding that age was positively associated with
742 children's English receptive vocabulary and productive morphology since we also observed that
743 older children tended to have higher English language scores.

744 There is an important finding regarding children with low language ability that must be
745 considered carefully. Although 42% of the children in this study were receiving speech/language
746 services at the time of data collection, all children generally tended to rate themselves as
747 speaking both languages well, although the degree of their ratings varied. It is crucial to design
748 questionnaires with multiple questions from a measurement perspective. For example, in looking
749 at the factor loadings (Table 3), the question "are you good at speaking Spanish?" was a strong
750 indicator of "self-perception of Spanish" (i.e., .93 loading), whereas this was slightly weaker for
751 "are you good at speaking English?" (i.e., .71 loading). These loadings can be roughly
752 interpreted similarly to correlations with the overall factor. Although children tended to respond
753 positively to both of these items, there was additional variation in their self-perceptions captured
754 by the other questionnaire items. From a questionnaire design perspective, we did not design the
755 Houston-Q to capture variation in language ability. We expected that even children with low
756 language ability (i.e., language disorders) would rate themselves as speaking well in at least one
757 of their languages. Our results suggested that this was the case. The Houston-Q cannot identify

758 children with low language ability but can potentially help identify a child, for example, with
759 stronger language proficiency in Spanish than in English and who has more experiences in
760 Spanish regardless of their language ability.

761

762 **Language Awareness**

763 The results of this study suggest that bilingual children have enough language awareness
764 to complete a questionnaire about their perceptions of their bilingual experiences and
765 proficiency. When children between the ages of 4 and 8 complete this questionnaire, they do so
766 reliably, and their responses are in general agreement with their proficiency in each language.
767 These results support previous studies that suggest that young children have enough language
768 awareness to self-report their relative language proficiency and bilingual experience (e.g.,
769 Babino & Stewart, 2016; Rojo & Echols, 2017). This finding is of interest because children are
770 usually not asked to provide this information, and instead, this information is often sought from
771 parents and teachers. We did not compare whether parents, teachers, or children provide the most
772 accurate information about the children, so we cannot make judgments about the overall
773 accuracy of the different reports. However, our results suggest that children might have a role in
774 providing this information because they are direct observers of their bilingual experience and
775 might be able to estimate their knowledge in each language compared to what parents and/or
776 teachers can report.

777

778 **Clinical Application**

779 An important piece of information during the assessment of language skills in bilingual
780 children is to understand how bilingual experience and proficiency in each language may play a

781 role in the child's overall language ability. This understanding is key to differentiating language
782 disorders from limitations or differences due to variability in proficiency and language exposure.
783 Administering the Houston-Q to children as part of the bilingual assessment could provide
784 important information about the child's perception of their current bilingual experience and
785 general proficiency in each language, which might facilitate identification of children's baseline
786 language experiences and strongest language prior to direct comprehensive language assessment.
787 Since this study included children with various levels of language ability, we recommend that
788 this questionnaire could be used by children with and without language disorders. Using the
789 child's self-reported bilingual experience and proficiency in each language may support clinical
790 assessment to consider the child's abilities in their strongest language for more accurate
791 identification of language disorders. However, it is important to note that this questionnaire was
792 not designed to identify children with low language ability.

793

794 **Limitations**

795 There are limitations to the interpretation of this study that are important to acknowledge.
796 This work provides preliminary evidence of the reliability and validity of the Houston-Q for
797 gaining some insight into Spanish-English speaking bilingual children's language experience and
798 proficiency. Although we believe the current scoring system is functional for clinical and
799 research use, further vetting with independent samples of bilingual children in the U.S. (and
800 other countries) is needed to better understand how children with different bilingual language
801 experiences respond to the Houston-Q. The questions may elicit different patterns of responses in
802 different contexts, and there may be outside factors that influence these patterns. For example,
803 question #10, which asks about having friends who speak Spanish and English, may be a more

804 effective indicator of self-perception of Spanish proficiency in areas with less bilingual language
805 support compared to Houston. Or, in contexts where Spanish is the primary societal language,
806 question #10 could reflect self-perception of English proficiency. These differences are essential
807 to examine carefully, to better understand the information that can be obtained from the
808 Houston-Q in various contexts.

809 Given the size of the current sample, we were not able to test for differences in scale
810 functioning by individual differences among children within the sample. Specifically, although
811 we examined overall associations between children's age and their subscale scores on the
812 Houston-Q, we did not have sufficient power to assess measurement invariance by language
813 ability level or age. Consequently, it is important to recognize that this study provides initial
814 evidence that children can complete the Houston-Q and that their responses broadly reflect
815 valuable information. Further specific examination of the scale (and subscale) functioning across
816 diverse samples of bilingual children, particularly among children at risk for language disorders,
817 is needed to inform the clinical utility of the measure in diagnostic contexts. For future users of
818 the Houston-Q, we recommend starting with the initial scoring system provided in this study. A
819 careful examination of the robustness of the provided item parameters will be needed to validate
820 it for use in other contexts and samples.

821 Finally, it is important to note the limitations of current measurement modeling,
822 particularly in quantifying distinct but related factors using a combination of dichotomous and
823 polytomous response options. We prioritized the establishment of a practical scoring approach
824 for the Houston-Q so that it could be easily used with basic computer software by both clinicians
825 and researchers. Specifically, we developed item weights for the Houston-Q are based on the
826 identified loadings from the categorical confirmatory factor analysis. The results of this work do

827 clearly suggest that this approach is preferable compared to weighting the items equally. Still, the
828 generalizability of the loadings is limited to the extent to which the participant sample is
829 representative. As more sophisticated techniques for scoring and representative sampling of
830 participants become more accessible, a more generalizable scoring approach may be
831 implemented to obtain scores quickly and reliably for individual children.

832

833

Conclusion

834 In this study, we examined the internal consistency reliability and preliminary criterion
835 validity of the Houston Questionnaire. The Houston-Q was created to gather information from
836 the child's perspective about their bilingual experience and proficiency in each language. Our
837 results provide evidence in support of the reliability and validity of the Houston-Q when used
838 with bilingual children between the ages of 4 and 8 with various levels of language ability and
839 different bilingual proficiency profiles.

840

841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871

References

- Adesope, O. O., Lavin, T., Thompson, T., & Ungerleider, C. (2010). A systematic review and meta-analysis of the cognitive correlates of bilingualism. *Review of Educational Research, 80*(2), 207–245. <https://doi.org/10.3102/0034654310368803>
- Archibald, L. M. D., & Joannisse, M. F. (2009). On the Sensitivity and Specificity of Nonword Repetition and Sentence Recall to Language and Memory Impairments in Children. *Journal of Speech, Language and Hearing Research, 52*(4), 899-914. [https://doi.org/10.1044/1092-4388\(2009/08-0099\)](https://doi.org/10.1044/1092-4388(2009/08-0099))
- Arias, G., & Friberg, J. (2017). Bilingual language assessment: Contemporary versus recommended practice in American schools. *Language, Speech, and Hearing Services in Schools, 48*(1), 1–15. https://doi.org/10.1044/2016_LSHSS-15-0090
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association
- Artiles, A. J., Rueda, R., Salazar, J., & Higareda, I. (2002). English- language learner representation in special education in California urban school districts. In D. Losen & G. Orfield (Eds.), *Racial inequality in special education* (pp. 117–136). Cambridge, MA: Harvard Education Press.
- Babino, A. & Stewart, M. A. (2016). “I like English better”: Latino dual language students’ investment in Spanish, English, and bilingualism. *Journal of Latinos and Education, 16*(1), 18–29. <https://doi.org/10.1080/15348431.2016.1179186>
- Bedore, L. M., & Peña, E. D., García, M., & Cortez, C. (2005). Conceptual versus monolingual scoring: When does it make a difference? *Language, Speech, and Hearing Services in Schools, 36*(3), 188-200. [https://doi.org/10.1044/0161-1461\(2005/020\)](https://doi.org/10.1044/0161-1461(2005/020))
- Bedore, L. M., & Peña, E. D. (2008). Assessment of bilingual children for identification of language impairment: Current findings and implications for practice. *International Journal of Bilingual Education and Bilingualism, 11*(1), 1–29. <https://doi.org/10.2167/beb392.0>
- Bedore, L. M., Peña, E. D., Gillam, R. B., & Ho, T.-H. (2010). Language sample measures and language ability in Spanish English bilingual kindergarteners. *Journal of Communication Disorders, 43*(6), 498–510. <https://doi.org/10.1016/j.jcomdis.2010.05.002>

872 Bedore, L. M., Peña, E. D., Joyner, D., & Macken, C. (2011). Parent and teacher rating of
873 bilingual language proficiency and language development concerns. *International Journal*
874 *of Bilingual Education and Bilingualism*, 14(5), 489–511.

875 Bedore, L. M., Peña, E. D., Summers, C. L., Boerger, K. M., Resendiz, M. D., Greene, K., ... &
876 Gillam, R. B. (2012). The measure matters: Language dominance profiles across
877 measures in Spanish–English bilingual children. *Bilingualism*, 15(3), 616–629.

878 Bishop, D. V., Snowling, M. J., Thompson, P. A., Greenhalgh, T., & Catalise Consortium.
879 (2016). CATALISE: A multinational and multidisciplinary Delphi consensus study.
880 Identifying language impairments in children. *PLOS One*, 11(7), e0158753.

881 Bridges, K., & Hoff, E. (2014). Older sibling influences on the language environment and
882 language development of toddlers in bilingual homes. *Applied Psycholinguistics*, 35(2),
883 225–241. <https://doi.org/10.1017/S0142716412000379>

884 Boateng, G. O., Neilands, T. B., Frongillo, E. A., Melgar-Quinones, H. R., & Young, S. L.
885 (2018) Best practices for developing and validating scales for health, social, and
886 behavioral research: A primer. *Frontiers in Public Health*, 6(149).
887 <https://doi.org/10.3389/fpubh.2018.00149>

888 Bohman, T. M., Bedore, L. M., Peña, E. D., Mendez-Perez, A., & Gillam, R. B. (2010). What
889 you hear and what you say: Language performance in Spanish-English bilinguals.
890 *International Journal of Bilingual Education and Bilingualism*, 13(3), 325–344.
891 <https://doi.org/10.1080/13670050903342019>

892 Castilla-Earls, A., Francis, D., Iglesias, A., & Davidson, K. (2019). The Impact of the Spanish-
893 to-English Proficiency Shift on the Grammaticality of English Learners. *Journal of*
894 *Speech, Language, and Hearing Research*, 62(6), 1739–1754.
895 https://doi.org/10.1044/2018_JSLHR-L-18-0324

896 Castilla-Earls, A., Bedore, L., Rojas, R., Fabiano-Smith, L., Pruitt-Lord, S., Restrepo, M. A., &
897 Peña, E. (2020). Beyond scores: Using converging evidence to determine speech and
898 language services eligibility for dual language learners. *American Journal of Speech-*
899 *Language Pathology*, 29, 1116–1132. https://doi.org/10.1044/2020_AJSLP-19-00179

900 Carroll, S. E. (2017). Exposure and input in bilingual development. *Bilingualism: Language and*
901 *Cognition*, 20(1), 3–16. <https://doi.org/10.1017/S1366728915000863>

902 de Ayala, R. J. (2013). *Factor analyses with categorical indicators*. In Y. Petscher, C.
903 Schatschneider, & D. L. Compton (Eds.), *Applied quantitative analyses in the educational*
904 *and social sciences* (pp. 208–242). New York, NY: Routledge.

905 De Houwer, A. (2014). *The absolute frequency of maternal input to bilingual and monolingual*
906 *children: A first comparison*. In T. Grüter & J. Paradis (Eds.), *Input and experience in*
907 *bilingual development* (pp. 37–58). Amsterdam, The Netherlands: John Benjamins.

908 De Houwer, A. (2017). Early multilingualism and language awareness. In J. Cenoz (Ed.),
909 *Encyclopedia of Language and Education* (pp. 83-94). Switzerland: Springer.

910 DeMars, C. E. (2012). Confirming testlet effects. *Applied Psychological Measurement*, 36(2),
911 104–121. <https://doi.org/10.1177/0146621612437403>

912 DiStefano, C., Zhu, M., & Mîndrilă, D. (2009). Understanding and using factor scores:
913 Considerations for the applied researcher. *Practical Assessment, Research and*
914 *Evaluation*, 14(20). <https://doi.org/10.7275/da8t-4g52>

915 Dunn, L. M., & Dunn, D. M. (2007). *PPVT-4: Peabody Picture Vocabulary Test-Fourth Edition*.
916 Pearson Assessments. <https://search.library.wisc.edu/catalog/999616587302121>

917 Dunn, L. M., Padilla, E. R., Lugo, D. E., & Dunn, L. M. (1986). *Tvip : Test De Vocabulario En*
918 *Imagenes Peabody : Adaptacion Hispanoamericana = Peabody Picture Vocabulary Test*
919 *: Hispanic-American Adaptation*. American Guidance Service.
920 <https://search.library.wisc.edu/catalog/999767172102121>

921 Dunn, T. J., Baguley, T., & Brunsten, V. (2014). From alpha to omega: A practical solution to
922 the pervasive problem of internal consistency estimation. *British Journal of Psychology*,
923 105(3), 399-412. <https://doi.org/10.1111/bjop.12046>

924 Duursma, E., Romero-Contreras, S., Szuber, A., Proctor, P., Snow, C., August, D., & Calderon,
925 M. (2007). The role of home literacy and language environment on bilinguals' English
926 and Spanish vocabulary development. *Applied Psycholinguistics*, 28(1), 171–190.
927 <https://doi.org/10.1017/S0142716406070093>

928 Farver, J. A. M., Lonigan, C. J., & Eppe, S. (2009). Effective early literacy skill development for
929 young Spanish-speaking English language learners: An experimental study of two
930 methods. *Child Development*, 80(3), 703–719. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-8624.2009.01292.x)
931 [8624.2009.01292.x](https://doi.org/10.1111/j.1467-8624.2009.01292.x)

932 Gollan, T. H., Weissberger, G. H., Runnqvist, E., Montoya, R. I., & Cera, C. M. (2012). Self-
933 ratings of spoken language dominance: A Multilingual Naming Test (MINT) and
934 preliminary norms for young and aging Spanish-English bilinguals. *Bilingualism:
935 language and cognition*, 15(3), 594-615.

936 Hammer, C. S., E. Komaroff, B. L. Rodriguez, L. M. Lopez, S. E. Scarpino, and B. Goldstein.
937 (2012) “Predicting Spanish–English Bilingual Children’s Language Abilities.” *Journal of
938 Speech, Language, and Hearing Research* 55 (5): 1251–1264.

939 Hoff, E., & Core, C. (2013). Input and language development in bilingually developing children.
940 *Seminars in Speech and Language*, 34(4), 215–226. [https://doi.org/10.1055/s-0033-
941 1353448](https://doi.org/10.1055/s-0033-1353448)

942 Hoff, E. (2017). How bilingual development is the same as and different from monolingual
943 development. *OLBI Working Papers*, 3–16.

944 Hoff, E., Burridge, A., Ribot, K. M., & Giguere, D. (2018). Language Specificity in the Relation
945 of Maternal Education to Bilingual Children’s Vocabulary Growth. *Developmental
946 Psychology*, 54(6), 1011-1019. <https://doi.org/10.1037/dev0000492>

947 Kan, P. F., & Windsor, J. (2010). Word learning in children with primary language
948 impairment: A meta-analysis. *Journal of Speech, Language, and Hearing Research*,
949 53(3), 739-756. [https://doi.org/10.1044/1092-4388\(2009/08-0248\)](https://doi.org/10.1044/1092-4388(2009/08-0248))

950 Kapantzoglou, M., Restrepo, M. A., Gray, S., & Thompson, M. S. (2015). Language ability
951 groups in bilingual children: A latent profile analysis. *Journal of Speech, Language, and
952 Hearing Research*, 58(5), 1549–1562. https://doi.org/10.1044/2015_JSLHR-L-14-0290

953 Kaufman, A. S., & Kaufman, N. L. (2004). *Kaufman Brief Intelligence Test* (Second Edition
954 ed.). Pearson, Inc.

955 Kohnert, K. (2010). Bilingual children with primary language impairment: Issues, evidence, and
956 implications for clinical actions. *Journal of Communication Disorders*, 43(6), 456–473.
957 <https://doi.org/10.1016/j.jcomdis.2010.02.002>

958 Leonard, L.B. (2014). *Children with specific language impairment.*, 2nd Edition. Cambridge,
959 MA: MIT Press.

960 Logan, J. A. R., Jiang, H., Helsabeck, N., & Yeomans-Maldonado, G. (2019, June 25). Should I
961 Allow my Confirmatory Factors to Correlate During Factor Extraction? Implications for
962 the Applied Researcher. <https://doi.org/10.31219/osf.io/zcsnv>

963 Lomax, R. G. (2013). *Introduction to structural equation modeling*. In Y. Petscher, C.
964 Schatschneider, & D. Compton (Eds.), *Applied quantitative analysis and the social*
965 *sciences* (pp. 245-264). New York: Routledge.

966 Lutz, A. (2008). Negotiating home language: Spanish maintenance and loss in Latino families.
967 *Latino(a) Research Review*, 6, 37-64.

968 Marian, V., Blumenfeld, H. K., & Kaushanskaya, M. (2007). The language experience and
969 proficiency questionnaire (LEAP-Q): Assessing language profiles in bilinguals and
970 multilinguals. *Journal of Speech, Language, and Hearing Research*, 50(4), 940–967.

971 Martin, B. (2012). Coloured language: identity perception of children in bilingual programmes.
972 <https://doi.org/10.1080/09658416.2011.639888>, 21(1–2), 33–56.
973 <https://doi.org/10.1080/09658416.2011.639888>

974 Melo-Pfeifer, S. (2015). Multilingual awareness and heritage language education: children’s
975 multimodal representations of their multilingualism.
976 <http://dx.doi.org/10.1080/09658416.2015.1072208>, 24(3), 197–215.
977 <https://doi.org/10.1080/09658416.2015.1072>

978 McNeish, D. (2017). Thanks coefficient alpha, we’ll take it from here. *Psychological Methods*,
979 23, 412–433. <http://dx.doi.org/10.1037/met0000144>

980 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavioral Research*
981 *Methods*, 22, 2287-2305. <https://doi.org/10.3758/s13428-020-01398-0>

982 Muthén, L. K., & Muthén, B. O. (1998-2019). *Mplus user’s guide* (8th Ed). Los Angeles, CA:
983 Múthen & Múthen

984 National Kids Count (2020). *Children living in linguistically isolated households by family*
985 *nativity in the United States*. The Annie E. Casey Foundation.
986 [https://datacenter.kidscount.org/data/tables/129-children-living-in-linguistically-isolated-](https://datacenter.kidscount.org/data/tables/129-children-living-in-linguistically-isolated-households-by-family-nativity)
987 [households-by-family-nativity](https://datacenter.kidscount.org/data/tables/129-children-living-in-linguistically-isolated-households-by-family-nativity)

988 Obied, V. M. (2009). How do siblings shape the language environment in bilingual families?
989 *International Journal of Bilingual Education and Bilingualism*, 12(6), 705–720.
990 <https://doi.org/10.1080/13670050802699485>

- 991 Osika, W., Friberg, P., & Wahrborg, P. (2007). A new short self-rating questionnaire to assess
992 stress in children. *International journal of behavioral medicine*, 14(2), 108–117.
993 <https://doi.org/10.1007/BF03004176>
- 994 Pratt, A. S., Peña, E. D., & Bedore, L. M. (2020). Sentence repetition with bilinguals with and
995 without DLD: Differential effects of memory, vocabulary, and exposure. *Bilingualism:
996 Language and Cognition*, 24(2), 305–318. <https://doi.org/10.1017/s1366728920000498>
- 997 Peña, E. D., Bedore, L. M., Gutierrez-Clellen, V. F., Iglesia, A., & Goldstein, B. A. (2008).
998 *Bilingual English-Spanish Assessment - Middle Extension Experimental Test Version
999 (Besa-Me)*. Unpublished manuscript.
- 1000 Peña, E. D., Bedore, L. M., & Kester, E. S. (2016). Assessment of language impairment in
1001 bilingual children using semantic tasks: two languages classify better than one.
1002 *International Journal of Language & Communication Disorders*, 51(2), 192–202.
1003 <https://doi.org/10.1111/1460-6984.12199>
- 1004 Peña, E. D., Bedore, L. M., Gutierrez-Clellen, V. F., Iglesia, A., & Goldstein, B. A. (2016).
1005 *Bilingual English-Spanish Assessment - Middle Extension Field Test Version (Besa-Me)*.
1006 Unpublished manuscript.
- 1007 Peña, E. D., Gillam, R. B., & Bedore, L. M. (2014). Dynamic assessment of narrative ability in
1008 English accurately identifies language impairment in English language learners. *Journal
1009 of Speech, Language, and Hearing Research*, 57, 2206–2220.
1010 https://doi.org/10.1044/2014_JSLHR-L-13-0151
- 1011 Peña, E. D., Gutierrez-Clellen, V. F., Iglesias, A., Goldstein, B., & Bedore, L. M. (2018).
1012 *Bilingual English-Spanish Assessment (Besa)*. Brookes Publishing.
- 1013 Perez-Leroux, A. T., Cuza, A., & Omas, D. (2011). From parental attitudes to input conditions
1014 Spanish-English bilingual development in Toronto. In K. Potowski (Ed.), *Bilingual
1015 youth: Spanish in English-speaking societies* (pp. 149–176). John Benjamins.
- 1016 Place, S., & Hoff, E. (2011). Properties of Dual Language Exposure That Influence 2-Year-Olds’
1017 Bilingual Proficiency. *Child Development*, 82(6), 1834–1849.
1018 <https://doi.org/10.1111/J.1467-8624.2011.01660.X>
- 1019 Restrepo, M. A. (1998). Identifiers of predominantly Spanish-speaking children with language
1020 impairment. *Journal of Speech, Language, and Hearing Research*, 41(6), 1398–1411.
- 1021 Revelle, W., & Condon, D. M. (2019). Reliability from alpha to omega: A tutorial.

1022 *Psychological Assessment*, 31(12), 1395 – 1411. <https://doi.org/10.1037/pas0000754>

1023 Rojas, R., Iglesias, A., Bunta, F., Goldstein, B., Goldenberg, C., & Reese, L. (2016). Interlocutor
1024 differential effects on the expressive language skills of Spanish-speaking English learners.
1025 *International journal of speech-language pathology*, 18(2), 166–177.
1026 <https://doi.org/10.3109/17549507.2015.1081290>

1027 Rojo, D. P., & Echols, C. H. (2017). Accepting labels in two languages: Relationships with
1028 exposure and language awareness. *OLBI Journal*, 8.
1029 <https://doi.org/10.18192/OLBIWP.V8I0.2115>

1030 Rujas, I., Mariscal, S., Murillo, E., & Lázaro, M. (2021). Sentence repetition tasks to detect and
1031 prevent language difficulties: A scoping review. *Children*, 8(7), 578.

1032 Samson, J., & Lesaux, N.K. (2009). Language minority learners in special education: Rates and
1033 predictors of identification for services. *Journal of Learning Disabilities*, 42(2), 148-1
1034 Semel, E., Wiig, E. H., & Secord, W. A. (2006). *Clinical Evaluation of Language*
1035 *Fundamentals–Fourth Edition, Spanish Version*. Pearson Education, Inc.

1036 Solans, M., Pane, S., Estrada, M. D., Serra-Sutton, V., Berra, S., Herdman, M., Alonso, J., &
1037 Rajmil, L. (2008). Health-related quality of life measurement in children and adolescents:
1038 a systematic review of generic and disease-specific instruments. *Value in health: the*
1039 *journal of the International Society for Pharmacoeconomics and Outcomes Research*,
1040 11(4), 742–764. <https://doi.org/10.1111/j.1524-4733.2007.00293.x>

1041 Strauss, M. E., & Smith, G. T. (2009). Construct validity: Advances in theory and methodology.
1042 *Annual Review of Clinical Psychology*, 5, 1-25.
1043 <https://doi.org/10.1146/annurev.clinpsy.032408.153639>

1044 Svalberg, A., M-L. (2007). Language awareness and language learning. *Language Teaching*,
1045 40(4), 287-308. <https://doi.org/10.1017/S0261444807004491>

1046 Tomoschuk, B., Ferreira, V. S., & Gollan, T. H. (2019). When a seven is not a seven: Self-
1047 ratings of bilingual language proficiency differ between and within language populations.
1048 *Bilingualism: Language and Cognition*, 22(3), 516-536.

1049 U.S. Census Bureau. (2021) 2019: American Community Survey 1-Year Subject Tables: S1601
1050 *Language Spoken at Home*. Retrieved from:
1051 [https://data.census.gov/cedsci/table?q=Houston&t=Language%20Spoken%20at%20Hom](https://data.census.gov/cedsci/table?q=Houston&t=Language%20Spoken%20at%20Home&tid=ACST1Y2019.S1601)
1052 [e&tid=ACST1Y2019.S1601](https://data.census.gov/cedsci/table?q=Houston&t=Language%20Spoken%20at%20Home&tid=ACST1Y2019.S1601)

- 1053 Vagh, S. B., Pan, B. A., & Mancilla-Martinez, J. (2009). Measuring growth in bilingual and
1054 monolingual children's English productive vocabulary development: The utility of
1055 combining parent and teacher report. *Child Development*, 80(5), 1545–1563.
1056 <https://doi.org/10.1111/j.1467-8624.2009.01350.x>
- 1057 Wiig, E. H., Semel, E., & Secord, W. A. (2013). *Clinical Evaluation of Language*
1058 *Fundamentals—Fifth Edition*. Pearson Education, Inc.
- 1059 Wood, Hoge, Schatschneider & Castilla-Earls (2018). Predictors of item accuracy on the *Test de*
1060 *Vocabulario en Imagenes Peabody* for Spanish-English speaking children in the United
1061 States. *International Journal of Bilingual Education and Bilingualism*,
1062 <https://doi.org/10.1080/13670050.2018.1547266>
- 1063
- 1064

1065 Table 1. Demographics and Language measure scores for children in the study ($n=113$)
 1066

	<i>n</i>	<i>M</i>	<i>SD</i>	%
Age (in months)	113	71.05	12.46	
Gender				
Male	64			56.6 %
Female	49			43.4 %
Mother's Level of Education				
No college	63			54.5 %
At least some college	50			45.5 %
Does the child qualify for free/reduced lunch?				
No	36			30.0 %
Yes	77			70.0 %
Child has received/is receiving services for speech/language?				
No	66			58.4 %
Yes	47			41.6 %
Language Spoken at Home				
English	12			10.6%
Spanish	56			49.6%
Both English and Spanish	45			39.8%
School Programs				
English-only	5			4.4%
Bilingual or Immersion	101			89.4%
Other: Saturday Spanish School	7			6.2%
Language Measures Norm-referenced assessments				
BESA/BESA-ME Morph Spanish		80.93	18.66	
BESA/BESA-ME Morph English		84.78	19.47	
BESA/BESA-ME Morph best language		92.19	17.06	
TVIP Spanish		86.72	17.93	
PPVT English		85.26	20.12	
CELF RO Spanish		6.75	3.11	
CELF SR English		6.76	3.56	

1067 *Note.* BESA/BESA-ME = Bilingual English-Spanish Assessment/Bilingual English-Spanish
 1068 Assessment-Middle Extension. Morph = Morphosyntax. TVIP = Test de Vocabulario en
 1069 Imágenes Peabody. PPVT = Peabody Picture Vocabulary Test, 4th Edition; CELF = Clinical
 1070 Evaluation of Language Fundamentals. RO = Recordando Oraciones / Recalling Sentences. SR =
 1071 Sentence Repetition.
 1072

Table 2

Fit Indices for Hypothesized Models Underlying Questionnaire

Model	χ^2	<i>df</i>	$\Delta\chi^2$	Δdf	$\Delta Sig.$	RMSEA	LB	UB	CFI	TLI
A 4-Factor with Items 10-11 crossed ¹	97.898	96		---		0.013	<.001	0.052	0.981	0.976
2-Factor with Items 10-11 crossed	103.470	101	5.895	5	.317	0.015	<.001	0.052	0.975	0.970
B 2-Factor with 10-11 on Spanish	105.984	103	2.625	2	.269	0.016	<.001	0.052	0.970	0.965
2-Factor with 11 only on Spanish ²	91.864	89		---		0.017	<.001	0.054	0.974	0.970
A 2-Factor: Bilingual Experience ¹	63.727	43		---		0.066	0.026	0.098	0.948	0.934
B 1-Factor: Bilingual Experience ²	64.360	44	0.233	1	.630	0.064	0.024	0.096	0.949	0.936

Note. $\Delta\chi^2$ is reported for the model comparisons against the previous (above) model.

¹Depicted in Figure 1A.

²Depicted in Figure 2. Finalized through discussion of item functioning, global fit, and consistency with theoretical expectations. The decrease in degrees of freedom reflects the full removal of question #10 from the measurement model.

Table 3

Standardized Item Loadings and Thresholds for Final Model

Factor	Questionnaire Item	Loading (SE)	Thresholds (SE)
Self-Perception of Spanish	1. Speak Spanish well (0/1)	0.93 (0.11)	-1.25 (0.16)
	2. Degree of speaking Spanish well	0.63 (0.10)	-1.13 (0.15)
			-0.84 (0.14)
			-0.60 (0.13)
			-0.49 (0.13)
	20. Spanish easiness (0/1)	0.79 (0.12)	-1.11 (0.15)
	21. Degree of Spanish easiness	0.69 (0.09)	-1.06 (0.15)
			-0.88 (0.14)
			-0.69 (0.13)
			-0.47 (0.12)
			-0.52 (0.12)
	6. Friends who speak Spanish (0/1)	0.37 (0.13)	-0.52 (0.12)
	7. Number of friends who speak Spanish	0.40 (0.14)	-0.77 (0.14)
			-0.28 (0.1)
			-0.04 (0.13)
			0.07 (0.13)
	24. Quantity of Spanish heard each day.	0.42 (0.12)	-0.96 (0.14)
			-0.62 (0.13)
			-0.32 (0.12)
			-0.16 (0.12)
	11. Number of Spanish-English speaking friends	0.27 (0.12)	-0.91 (0.14)
			-0.65 (0.13)
			-0.11 (0.12)
			0.23 (0.12)
Factor	Questionnaire Item	Loading (SE)	Thresholds (SE)
Self-Perception of English	3. Speak English well (0/1)	0.71 (0.16)	-1.35 (0.17)
	4. Degree of speaking English well	0.60 (0.12)	-1.33 (0.17)
			-0.93 (0.14)
			-0.70 (0.13)
			-0.43 (0.13)
	22. English easiness (0/1)	0.58 (0.14)	-0.76 (0.13)
	23. Degree of English easiness	0.45 (0.19)	-1.02 (0.15)
			-0.77 (0.14)
			-0.34 (0.13)
			-0.05 (0.12)

	8. Friends who speak English (0/1)	0.56 (0.17)	-0.77 (0.13)
	9. Number of friends who speak English	0.17 (0.16)	-0.74 (0.14)
			-0.32 (0.13)
			-0.18 (0.13)
			-0.08 (0.13)
	25. Quantity of English heard each day.	0.45 (0.13)	-1.11 (0.15)
			-0.62 (0.13)
			-0.223 (0.12)
			0.05 (0.12)
Factor	Questionnaire Item	Loading (SE)	Thresholds (SE)
Bilingual Experience	5a. Language used with a family member (1).	0.80 (0.06)	-0.18 (0.12)
			1.11 (0.15)
	5b. Language used with a family member (2).	0.42 (0.10)	-0.34 (0.12)
			0.76 (0.13)
	5c. Language used with a family member (3).	0.74 (0.07)	-0.01 (0.12)
			0.88 (0.14)
	12. Language spoken with bilingual friends.	0.65 (0.07)	-0.45 (0.12)
			0.45 (0.12)
	13. Language used with teacher.	0.24 (0.11)	-0.82 (0.14)
			0.88 (0.14)
	14. Language used for learning to write.	0.56 (0.08)	-0.79 (0.13)
			1.15 (0.15)
	15. Language used for watching TV.	0.62 (0.08)	-1.01 (0.15)
		0.52 (0.13)	
16. Language used when playing at the park.	0.79 (0.06)	-0.56 (0.13)	
		0.51 (0.13)	
17. Language used at parties/family reunions.	0.65 (0.08)	-0.41 (0.13)	
		0.66 (0.13)	
18. Language used to read books.	0.62 (0.07)	-0.51 (0.13)	
		1.01 (0.15)	
19. Language used for learning to read.	0.52 (0.09)	-0.58 (0.13)	
		1.12 (0.15)	

Note. The underlying latent trait mean was set to zero, with a variance of 1.

For the Bilingual Experience latent factor, -1 = Experience in Spanish, 0 = Experience in Spanish and English, and 1 = Experience in English.

Table 4

Means, standard deviations, and correlations with confidence intervals

Variable	<i>M</i>	<i>SD</i>	1	2	3	4	5	6	7	8	9	10
Houston-Q	1. Spanish SP	7.73	2.15									
	2. English SP	7.69	2.09	-.24*								
	3. Bilingual Exp	8.94	4.53	-.61**	.42**							
Spanish	4. CELF RO	6.75	3.11	.36**	-.10	-.27**						
	5. TVIP	86.72	17.93	.23*	-.08	-.27**	.64**					
	6. BESA Morph	80.93	18.66	.42**	-.11	-.35**	.81**	.69**				
English	7. CELF SR	6.76	3.56	-.43**	.32**	.36**	.30**	.19*	.16			
	8. PPVT	85.26	20.12	-.40**	.24*	.38**	.07	.21*	.05	.73**		
	9. BESA Morph	84.78	19.47	-.39**	.23*	.35**	.12	.15	.11	.78**	.80**	
10. Best BESA	92.19	17.06	-.22*	.15	.14	.44**	.40**	.50**	.70**	.61**	.78**	
11. Age (mos)	70.05	12.46	-.19*	-.13	.02	-.15	.01	-.05	.10	.21*	.47**	.38**

Note. *M* and *SD* are used to represent mean and standard deviation, respectively. SP = Self-Perception. Exp = Experience. CELF RO = Clinical Evaluation of Language Fundamentals. RO = Recordando Oraciones / Recalling Sentences. TVIP = Test de vocabulario en imagenes Peabody. BESA = Bilingual English-Spanish Assessment. Morph = Morphosyntax. SR = Sentence Repetition. PPVT = Peabody Picture Vocabulary Test, 4th Edition. * indicates $p < .05$. ** indicates $p < .01$