

Measuring a Basic Developmental Vocabulary Across Many Languages

Anonymous CogSci submission

Abstract

Early language skill is predictive of later life outcomes, and is thus of great interest to developmental psychologists and clinicians. The Communicative Development Inventories (CDIs), including a parent-reported inventory of early-learned vocabulary items, has proven to be a valid and reliable instrument for measuring children's early language skill. CDIs have been adapted to dozens of languages, and cross-linguistic comparisons thus far show both consistency and variability in language acquisition trajectories. Here, we use item-response theory models trained on 26 languages to examine the psychometric properties of translation-equivalent concepts that are frequently included on CDIs, with the goal of identifying a short list of concepts that are of similar learning difficulty in the majority of the languages, in order to propose a pool list of 'universal' words that can be used as a starting point for future CDI adaptations. After identifying a list of 96 items, we test how well this list generalizes to an additional 8 languages.

Keywords: early language learning; CDI; psychometrics; cross-linguistic comparison; Swadesh vocabulary;

Introduction

Tools that enable valid assessments of children's early language abilities are invaluable for researchers, clinicians and parents, as early language skill is predictive of educational outcomes years later (e.g., Bleses, Makransky, Dale, Højen, & Ari, 2016). The MacArthur-Bates Communicative Development Inventories [CDIs; Fenson et al. (2007); (marchman2023?)] are parent report assessments that provide reliable and valid estimates of children's early vocabulary size and other aspects of early communicative development, such as gesture use and early use of word combinations. Parent report is a relatively quick and low-cost method to assess early language skills since it takes advantage of the fact that parents are "natural observers" of their child's skills and does not depend on a child engaging with an (unfamiliar) experimenter. Over the years, the CDIs have been adapted to dozens of languages, with forms now available in English, Spanish, French, Hebrew, and Mandarin, to name just a few. Recently, data from more than 85,000 CDIs in 38 languages have been archived in a central repository [Wordbank; Frank, Braginsky, Yurovsky, & Marchman (2017)]. These data have revealed both cross-linguistic consistency and variability in early language skills, with insights from these patterns informing theories of early language learning (Frank, Braginsky, Yurovsky, & Marchman, 2021). For example, cross-linguistic analyses indicate that measures of vocabulary size

are tightly correlated with other aspects of early language skill, like gesture and grammatical competence. Thus, over development, the language system is "tightly woven" (Bates et al., 1994; Frank et al., 2021) and early vocabulary size serves as a good proxy measure of children's overall language skill.

On the CDIs, vocabulary size is assessed via a checklist format, which enables caregivers to scan and recognize words their child produces or understands, rather than relying on recall alone. For example, the American English CDI Words & Sentences (CDI:WS) form, targeting children 16-30 months of age, is comprised of 680 words from 22 semantic categories, including nouns (e.g., Body Parts, Toys, and Clothing), action words (e.g. verbs), descriptive words (e.g. adjectives), and closed-class words such as pronouns. The vocabulary checklist from the American English CDI Words & Gestures (CDI:WG) form, targeting younger children ages 8 to 18 months, is comprised of ~400 words from a similar set of categories. Items on these original forms were chosen to reflect a range of difficulty levels (i.e., both easy, moderate, and more difficult), as well as capture the linguistic and societal contexts of (most) children living in the US.

Short versions of each of these forms are also available, each with about 100 items (e.g., Fenson et al., 2000), consisting of a set of items that generate scores that more strongly correlate with scores on the long forms, while retaining representation across a broad set of semantic categories.

Following the guidelines¹ from the MacArthur-Bates CDI Advisory Board, the process of adapting a CDI for a language other than American English goes well beyond simply translating items on these forms to that new language. While the process can begin with identifying translation equivalents (i.e., items that capture the same general concept in both languages, e.g., "dog" in English, and "perro" in Spanish), the final item set must then be filtered so that all items appropriately reflect the linguistic and sociocultural context of the children learning that language. This process usually requires considerable time and effort by researchers who are native speakers of a language, to first select and identify translation equivalents and to then iteratively add, refine, and pilot the new CDI in the target language. Because the goal is to obtain the set of items that best capture general trends and individual

¹<http://wordbank.stanford.edu/contributors>

differences in that language, the items across CDIs in different languages do not necessarily overlap to a great extent. For example, the American English CDI:WS and Mexican Spanish CDI:WS forms – two of the first CDIs created – each have 680 words, but only have 463 overlapping concepts (68%).

Given that it is well-established that, all over the world, early learned words reflect the people and things that children are likely to experience, that is, words for family members, animals, and common household objects (Frank et al., 2021), it is reasonable to ask: Is there a single set of translation equivalents that would meet the criteria for inclusion on CDIs from multiple languages? To facilitate this effort, it is useful to leverage recent work in Item-Response Theory [IRT; Embretson & Reise (2013)] models. IRT models infer both the abilities of test takers and the difficulty of individual test items (i.e., words), along standardized dimensions. Recent work using IRT models have facilitated our understanding of the psychometric properties of CDIs instruments. As such, they offer the potential to not only yield more accurate measures of children’s language ability, but also to enable the construction of Computerized Adaptive Tests (CATs), which choose the next test item based on the responses to the previous items, and thus quickly hone in on the test-takers language ability. CAT-based CDIs presenting 50 or fewer items have been found to strongly correlate with scores on the full CDI:WS (Chai, Lo, & Mayor, 2020; Makransky, Dale, Havmose, & Bleses, 2016; Mayor & Mani, 2019). A general method for creating CDI CATs that work well across a broader age range (12-36 months) has been proposed, and tested for American English and Mexican Spanish (Kachergis, Marchman, Dale, Mankewitz, & Frank, 2022). However, the IRT model driving each CAT needs to be trained on a large and normative dataset, which may not be available in a given language. To date, the IRT models are also fitted separately for each language, and the fitted parameters (e.g., word difficulty) are likely to vary across languages.

The goal of the current study is to use IRT modeling in conjunction with data from Wordbank to examine whether there might be a core set of concepts that are frequently included on CDIs, and—importantly—whether many of them are of roughly equal difficulty across many languages. This work takes its inspiration from the fields of lexicostatistics and glottochronology, where researchers (notably, 1971?) have proposed a list of ‘universal concepts’—concepts that exist in all catalogued languages—in order to quantify the genealogical relatedness and dates of divergence of languages. For example, the original Swadesh list contains 100 words, comprised of categories including common pronouns (“I”, “you”, “we”), animals (“man”, “fish”, “bird”, “dog”), objects (“tree”, “leaf”, “sun”, “mountain”), and verbs (“kill”, “die”, “see”, “sleep”). Extending this work to the development of a universal CDI, or “Swadesh CDI,” would include all of the concepts that researchers have chosen to include on the vast majority of CDI:WS adaptations, and which have relatively similar difficulty across many languages. More broadly, this

work examines which types of words (and their corresponding concepts) are more or less similar across languages in terms of the ease with which they are learned, revealing commonalities as well as idiosyncrasies in children’s early experiences. We can ask: Which semantic categories, e.g., animals, household objects, food and drink, or another category, are most consistently learned across languages? Which are more variable? These types of cross-linguistic comparisons may give new insight into theoretical questions surrounding the similarity and differences between language experience and development in different cultural and linguistic contexts.

In particular, our contributions are 1) to fit IRT models to 28 CDI:WS datasets, 2) to identify 96 candidate “Swadesh CDI” items from a cross-linguistic comparison of concept difficulty and inclusion, and 3) to test how well this Swadesh CDI generalizes to a set of 8 additional low-data languages. We end by making a concrete proposal for how this Swadesh CDI list could be used in creating future CDI adaptations, and by discussing the strengths and weaknesses of our approach.

Methods

Item Response Theory

A variety of IRT models targeting different types of testing scenarios have been proposed (see Baker, 2001 for an overview), but for the dichotomous responses that parents make for each item (word) regarding whether their child can produce a given word, we will use the popular 2-parameter logistic (2PL) model that we have previously found is best justified for CDI data (see Kachergis et al., 2022).

The 2PL model jointly estimates for each child j a latent ability θ_j (here, language skill), and for each item i two parameters: the item’s difficulty b_i and discrimination a_i , described below. In the 2PL model, the probability of child j producing a given item i is

$$P_i(x_i = 1 | b_i, a_i, \theta_j) = \frac{1}{1 + e^{-Da_i(\theta_j - b_i)}}$$

where D is a scaling parameter ($D = 1.702$) which makes the logistic more closely match the ogive function used in a standard factor analysis (Chalmers, 2012; Reckase, 2009). Children with high latent ability (θ) will be more likely to produce any given item than children with lower latent ability, and more difficult items will be produced by fewer children (at any given θ) than easier items. The discrimination (a_i) adjusts the slope of the logistic (in the classic 1-parameter logistic (1PL or Rasch) model, the slope is always 1). Items with higher discrimination (i.e. slopes) better distinguish children above vs. below that item’s difficulty level, and hence are generally more useful. While other standard IRT models exist (e.g., the 3-parameter logistic model adds a ‘guessing’ parameter for each test item), we have found elsewhere (Kachergis et al., 2022) that the 2PL model is best supported by the Wordbank data.

Datasets

Language	WS items	WS N
Norwegian	731	9304
English (American)	680	8828
Danish	725	3714
Portuguese (European)	639	3012
Turkish	711	2422
Mandarin (Taiwanese)	696	1897
Spanish (Mexican)	680	1853
English (Australian)	558	1520
Korean	641	1376
Cantonese	804	1295
German	588	1181
Slovak	609	1066
Mandarin (Beijing)	799	1056
Russian	728	1037
French (Quebecois)	664	929
Swedish	710	900
Spanish (Argentinian)	699	784
Italian	670	752
French (French)	690	665
Spanish (European)	588	593
Hebrew	605	518
Latvian	723	500
Czech	553	493
Croatian	717	377
Hungarian	802	363
Dutch	704	303
Greek (Cypriot)	815	176
Spanish (Peruvian)	600	105
Kigirama	696	100
English (Irish)	660	99
Irish	691	99
Kiswahili	705	90
Finnish	581	70
Persian	558	50

Table 1: Number of CDI:WG and CDI:WS items and subjects (N) per language.

We report IRT analyses for twenty-six languages from Wordbank (Frank et al., 2017), comprising production data from CDI:WS vocabulary checklists.² Data from the first twenty-six rows of Table 1 (Norwegian through Dutch) will be used to select a pool of words with approximately equal cross-linguistic difficulty. CDI:WS production data from an additional eight languages (Table 1: Greek (Cypriot) through Persian) had too few participants to be analyzed with IRT, but will be used to test how well the selected pool of generalize to new languages.

Participants The CDI:WS production dataset consists of the combined Wordbank production data for 47527 children

aged 16-30 months on 23020 items across 34 forms. Figure 1C and 1D shows children’s production scores vs. age for the CDI:WG and CDI:WS datasets respectively. Note that the distributions of demographic variables (age, sex, maternal education, etc.) of these datasets are not matched, so comparing overall language ability estimates across languages would be impossible. (See Frank et al. (2021) for a discussion of effects of demographic variables on vocabulary development.) Thus, we will focus only on the estimated item parameters, and in particular the variability of item difficulty (b_i).

Instruments When a CDI:WS forms was administered, caregivers were asked to indicate for each vocabulary item on the instrument whether or not their child can recognizably produce (say) the given word.

“Produces” responses were coded as 1 and all other responses were coded as 0. Our datasets consist of a dichotomous-valued response matrix for each language, of size N subjects \times W words.

Results

All models, data, code for reproducing this paper are available on OSF³. Across the 26 IRT models for different CDI:WS forms, parameters for a total of 17715 items were obtained. Figure 1 shows the average difficulty of items in each language, but as mentioned above these differences should not be interpreted, as they may be related to differences in the samples. However, it is useful to show them to emphasize once more that our main focus will be on the cross-linguistic *variability* in words’ difficulty (which should be immune to random effects from sample differences across languages), rather than differences in mean difficulties.

Figure 2 shows the average cross-linguistic difficulty of CDI items by semantic category. Sounds (e.g., animal noises), body parts, and common nouns tend to be early-learned, and thus have higher easiness values, while abstract words such as time words and morphologically complex helping verbs are later-learned (i.e., difficult).

Cross-linguistic similarities

We look at the Spearman correlation between the item difficulty of each language compared to each other language. We might expect this to recapitulate the historical relationship between languages, with more similar languages having more similar item difficulties (e.g., Quebecois and European French). (cut/move to appendix?)

Identifying Swadesh CDI Candidates CDI:WS forms have a median of 693 uni-lemmas defined (range: 553 (Czech) to 804 (Cantonese)), and there are a total of 1840 uni-lemmas defined across the 26 languages. Only 61 uni-lemmas appear on all 28 CDI:WS forms: too few to comprise a short form. To expand the pool of items, we consider thresholds on both 1) the variability of an item’s difficulty across languages (lower is better), and 2) the number of CDI:WS forms

²<http://wordbank.stanford.edu/contributors>

³OSF repository: <https://osf.io/XXX/>.

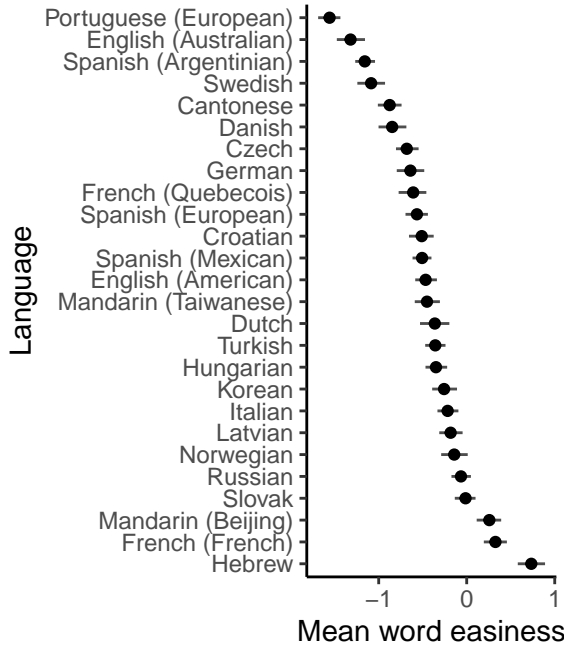
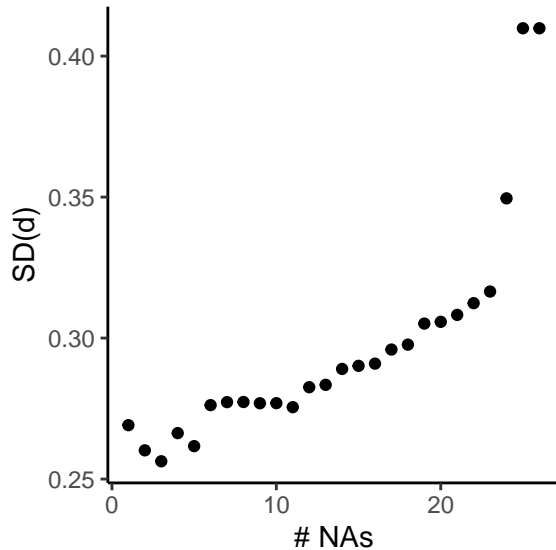


Figure 1: Mean word difficulty per language. Bars represent bootstrapped 95% confidence intervals.

on which it appears (more is better). Below we examine the standard deviation of uni-lemmas' cross-linguistic difficulty as a function of how many languages that uni-lemma is missing from. For now we consider items with less than median variability in difficulty that are included on 10 or more of the 26 CDI:WS forms.



Of these, a total of 1312 uni-lemmas are included in more than one language, and only 656 uni-lemmas are included in 10 or more of the languages. We will start by considering this more restricted list.

To evaluate how variable items are in their cross-linguistic difficulty, we calculate the standard deviation (SD) of each

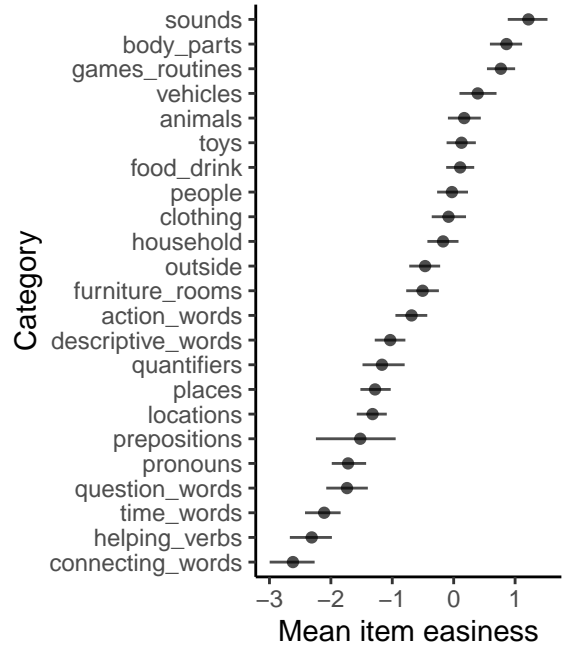


Figure 2: Mean cross-linguistic difficulty of CDI words by semantic category. Bars represent bootstrapped 95% confidence intervals.

uni-lemma's difficulty. The median SD is 1.17 ($SD=0.41$), so we consider the items with SD less than half the median SD (i.e., $SD < 0.96$). These 106 candidate "Swadesh CDI" items are summarized by semantic category below (see OSF for full list). 52% of Swadesh CDI words are nouns, 21% are predicates, 22% are function words, and 5% are other. This breakdown is comparable to the lexical category percentages on the 680-item English CDI:WS (46% nouns, 24% predicates, 15% function words, and 15%), minimally suggesting that this method of selecting candidate items is not biased by lexical category. The candidate items, ultimately selected by lower variability, were also present on more forms than typical in the selection set: on average, each item appeared on 20.0754717 forms.

Comparing the IRT parameters of the Swadesh CDI words to the rest of the items showed that the candidate Swadesh items are significantly easier (mean Swadesh $d = 0.03$, other items' mean $d = -0.51$, $t(2737) = 12.74$, $p < .001$). If the Swadesh list is too focused on early-learned items, this may result in ceiling effects for older children. It may be prudent to consider expanding the Swadesh list to include some more difficult items. The discrimination parameter (i.e., slope) of the Swadesh items did not significantly differ from the other items.

Comparing Swadesh CDI to Full CDI:WS

How well do sumscores from the Swadesh CDI items correlate with full CDI:WS scores? On average, for the 26 languages the IRT analysis was based on, the Swadesh CDI's

sumscores were strongly related to the full CDI:WS scores (mean $r = 0.990$; $min = 0.976$, $max = 0.996$). However, these correlations were slightly but significantly lower than scores based on randomly-selected items (mean $r = 0.993$, $min = 0.990$, $max = 0.996$, paired $t(25) = 5.67$, $p < .001$). This is likely due to some ceiling effects for older children on the Swadesh CDI, since the randomly-selected items are likely somewhat easier on average than the Swadesh items (see above results).

For the eight low-data languages, a similar comparison revealed that the Swadesh CDI's sumscores were again strongly related to the full CDI:WS scores (mean $r = 0.978$; $min = 0.959$, $max = 0.992$). As with the training set, however, the Swadesh list correlations were slightly but significantly lower than scores based on randomly-selected items (mean $r = 0.985$; $min = 0.972$, $max = 0.992$; $t(7) = 3.22$, $p = .01$).

(Maybe show plot of Swadesh vs. full CDI scores, by age? color by language? might be overwhelming..)

ToDo: use average IRT parameter per item to calculate ability rather than sumscore? propose a way of picking more difficult items? or maybe the value here is in finding the universal (easy) words, and it is up to expert native speakers to choose the more difficult (and idiosyncratic) words for each language.

If I had to find another way of getting the best cross-linguistic words, I would 1) run simulated CATs in all of the languages, and 2) pick the words that cropped up the most often across all languages.

Discussion

Acknowledgements

We would like to thank all of the contributors to Wordbank, from the researchers who created and adapted the CDIs to those who collected the data (as well as the participants), to those who have created and maintained Wordbank over the years.

References

- 10 Baker, F. B. (2001). *The basics of item response theory*. ERIC.
- Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P. S., Reznick, J. S., ... Hartung, J. (1994). Developmental and stylistic variation in the composition of early vocabulary. *Journal of Child Language*, 21(1), 85–123.
- Bleses, D., Makransky, G., Dale, P. S., Højen, A., & Ari, B. A. (2016). Early productive vocabulary predicts academic achievement 10 years later. *Applied Psycholinguistics*, 37(6), 1461–1476.
- Chai, J. H., Lo, C. H., & Mayor, J. (2020). A Bayesian-inspired item response theory-based framework to produce very short versions of MacArthur-Bates Communicative Development Inventories. *Journal of Speech, Language, and Hearing Research*, 63(10), 3488–3500.
- Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, 48(6), 1–29. <http://doi.org/10.18637/jss.v048.i06>
- Embretson, S. E., & Reise, S. P. (2013). *Item response theory*. Psychology Press.
- Fenson, L., Marchman, V. A., Thal, D. J., Dale, P. S., Reznick, J. S., & Bates, E. (2007). *MacArthur-Bates Communicative Development Inventories: User's guide and technical manual (2nd ed.)*. Baltimore, MD: Brookes.
- Fenson, L., Pethick, S., Renda, C., Cox, J. L., Dale, P. S., & Reznick, J. S. (2000). Short-form versions of the MacArthur Communicative Development Inventories, 21, 95–116.
- Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank: An open repository for developmental vocabulary data. *Journal of Child Language*, 44(3), 677.
- Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V.

Actions	dance, sing, blow, draw, swim, wait, catch, tickle, lick, stay, knock
Animals	cat, cow, <i>dog</i> , lion, animal, ant, chicken, spider, goat, turkey
Body Parts	<i>nose</i> , <i>tongue</i> , <i>tooth</i> , head, lip, shoulder, penis, vagina
Clothing	button, necklace, belt, skirt
Connectives	and
Descriptives	blue, green, yellow, black, first, heavy, sticky, dirty, awake
Food/Drink	milk, yogurt, peas, potato chips, lollipop, hamburger
Furniture/ Rooms	Door, washing machine, TV, drawer, garage
Games/ Routines	yes, shh, breakfast, shopping
Household	key, spoon, broom, hammer, radio, toothbrush, paper, glasses, nail, pacifier
Locations	in, under, behind, back
Outside	moon, <i>star</i> , sky, hose
People	mommy, daddy, doctor, pet's name, child's name, police, mailman
Places	school, church, gas station
Sounds	meow, cockadoodledoo, grrr, choo choo
Time Words	yesterday, tonight, tomorrow
Vehicles	airplane, helicopter, firetruck

Figure 3: The 96 proposed Swadesh CDI words by semantic category.

- A. (2021). *Variability and consistency in early language learning: The wordbank project*. MIT Press.
- Kachergis, G., Marchman, V. A., Dale, P. S., Mankewitz, J., & Frank, M. C. (2022). Online computerized adaptive tests of children's vocabulary development in english and mexican spanish. *Journal of Speech, Language, and Hearing Research*, 65(6), 2288–2308.
- Makransky, G., Dale, P. S., Havmose, P., & Bleses, D. (2016). An item response theory-based, computerized adaptive testing version of the MacArthur-Bates Communicative Development Inventory: Words & Sentences (CDI:WS). *Journal of Speech, Language, and Hearing Research*, 59(2), 281–289.
- Mayor, J., & Mani, N. (2019). A short version of the MacArthur-Bates Communicative Development Inventories with high validity. *Behavior Research Methods*, 51(5), 2248–2255.
- Reckase, M. D. (2009). Multidimensional item response theory models. In *Multidimensional item response theory* (pp. 79–112). Springer.