

**Wordbank: An Open Repository for Developmental Vocabulary Data**

### **Abstract**

Understanding the processes by which vocabulary grows provides a window into linguistic and cognitive development. The MacArthur-Bates Communicative Development Inventories (CDIs) are a widely-used family of parent-report instruments for easy and inexpensive data-gathering about early language acquisition. CDI data have been used to explore a variety of theoretically-important topics, but with few exceptions, researchers have had to rely on data collected in their own lab: there is no resource that offers researchers the opportunity to share and access raw CDI data. We describe Wordbank, a structured database of vocabulary measures from CDI forms combined with a browsable web interface. Wordbank archives CDI data across languages and labs, providing a resource for researchers interested in early language as well as a platform for novel analyses.

## Introduction

Learning language is one of the most impressive and intriguing human accomplishments: A child’s first words are eagerly awaited by parents, dutifully recorded in baby books, and celebrated with family and friends. Early achievements in vocabulary are correlated with subsequent progress in grammar (Bates & Goodman, 1999), and understanding the processes by which vocabulary grows provides a window into mechanisms of linguistic and cognitive development more generally (Bloom, 2002). In addition, since vocabulary growth lays the groundwork for future academic achievement (Hart & Risley, 1995; Marchman & Fernald, 2008), a better scientific understanding of factors affecting vocabulary development holds the promise of societal benefits.

The MacArthur-Bates Communicative Development Inventories (CDIs; Fenson et al., 1994, 2007) are a widely-used family of parent-report instruments for easy and inexpensive data-gathering about early language acquisition. CDI data have been used to explore many theoretically-rich topics, including variation in early word production (Fenson et al., 1994), vocabulary composition (Bates et al., 1994), and the growth of lexical networks (Hills, Maouene, Maouene, Sheya, & Smith, 2009). With few exceptions, however, researchers have had to rely on data collected in their own lab. While CDI norms are available (Fenson et al., 2007; Jørgensen, Dale, Bleses, & Fenson, 2010), no public resource offers researchers the opportunity to share and access raw, cross-linguistic data at the scale necessary to address questions about demographic variation, vocabulary composition, relations with grammatical development, and other important issues.

To remedy this issue, we introduce Wordbank ([wordbank.stanford.edu](http://wordbank.stanford.edu)), a structured database of developmental vocabulary data. Building on previous tools like CLEX (Jørgensen et al., 2010), Wordbank archives raw CDI data across languages and labs, providing a large-scale database of information about children’s vocabulary knowledge. The site hosts an interactive and expandable set of in-depth analyses that can

be explored by interested researchers, students, and members of the public. Wordbank lowers the cost of new, exploratory analyses by facilitating the productive reuse of data.

The current paper presents the Wordbank site in detail. We begin by discussing the motivations for constructing such a site. The bulk of the paper then describes the Wordbank site, including its database architecture, its web-based front-end, and its extensibility. We end by presenting `wordbankr`, a package for the R statistical programming language that allows research users to access the database directly, showing as a case study an analysis of gender differences in productive vocabulary across languages.

## Motivation and Background

The nature and course of early word learning is a window into children’s growing understanding of the world around them. Infants typically begin to associate sound sequences with meanings toward the end of their first year of life, responding appropriately to familiar phrases such as “It’s time for bath!” (e.g., by plopping down and attempting to remove their shoes); the production of recognizable words begins somewhat later. Early words cross-cut a variety of linguistic categories, but generally consist of names for caregivers (e.g., mama), common objects (e.g., bottle, shoe), social expressions (e.g., bye-bye), and actions or routines (e.g., peekaboo, throw) (Nelson, 1973; Tardif et al., 2008).

New words enter children’s expressive vocabularies over the next several months, starting slowly but quickly accelerating. Children reach an average of 300 words by 24 months and more than 60,000 by the time they graduate from high school (Fenson et al., 2007). At the same time, there are significant individual differences in language acquisition. For example, according to detailed observational studies, although some 18-month-olds already produce 50–75 words, others produce no words at all, and will not do so until they are 22 months or older (e.g., Brown, 1973). Does this

variability have downstream effects on later development? And how can such differences be measured accurately and efficiently?

### *Measuring early vocabulary*

Traditional studies of language development typically apply a combination of observational assessment and structured tests, frequently relying on short samples of interactions and small samples of children. Discerning both the universal features and natural variation of language development has been greatly facilitated by the development of parent report instruments like the MacArthur-Bates CDI (Fenson et al., 1994, 2007) and the Language Development Survey (Rescorla, 1989). The CDIs were developed across a period of more than 40 years. Originally designed for use in a research study (Bates, 1976), the instruments have evolved from a structured interview to the current paper-and-pencil format and are now increasingly administered online (e.g., Kristoffersen et al., 2013).

While other assessment tools exist for slightly older children (e.g., the Peabody Picture Vocabulary Test 4; Dunn & Dunn, 2007), to our knowledge, no other measure allows cost-effective global language assessment for children in the critical age ranges between the emergence of language and the period when children become more able to engage in structured, face-to-face activities (around 30 months). Naturalistic observations are the leading candidate, but such observations are extremely costly and time-consuming to transcribe and annotate, difficulties reflected in the small sample sizes for many studies using this method. In addition, naturalistic observations do not measure children's language comprehension, a variable of interest for many early language researchers. Estimates of productive vocabulary from naturalistic observation are highly correlated with the CDI within studies (e.g. Bornstein & Haynes, 1998), but affected substantially by length of the session, context, and interlocutor when comparing across studies.

Parent report measures like the CDI and LDS take advantage of the fact that parents are expert observers of their child. CDI instruments ask about use of communicative gestures, grammar, and symbolic play, as well as vocabulary, which is measured using checklists consisting of representative samples of words. Parents choose the words their child currently “understands” (comprehension, measured for younger children) or “says” (production, measured for both younger and older children). The checklists contain words from many different semantic (e.g., animals names, household items) and syntactic (e.g., action words, connectives) categories, resulting in broader samples of lexical knowledge than are available from other methods. The instruments come in two versions: Words & Gestures (8–18 months) and Words & Sentences (16–30 months). Originally designed for English, parallel instruments have now been adapted for more than 60 languages.

#### *Limitations of parent report*

Although the standardized parent report format contributes to the availability of large amounts of data in a comparable format, there are limitations to the parent report methodology as well (Tomasello & Mervis, 1994; Feldman et al., 2000). First, parents may be biased observers; some may overestimate, while others likely underestimate their children’s abilities. There is some evidence that some variability may be due to reporting biases linked to factors such as SES (Feldman et al., 2000; Fenson et al., 2000; Feldman et al., 2005). Second, parent reports of comprehension for younger children likely suffer from a number of biases and are probably substantially more accurate for content words than function words. Third, the items on the original CDI instruments were chosen to be a representative sample of vocabulary for the appropriate age and language (Fenson et al., 1994), not with the intention that they would be a complete set of words that could be compared across instruments or that they would be individually reliable and license the

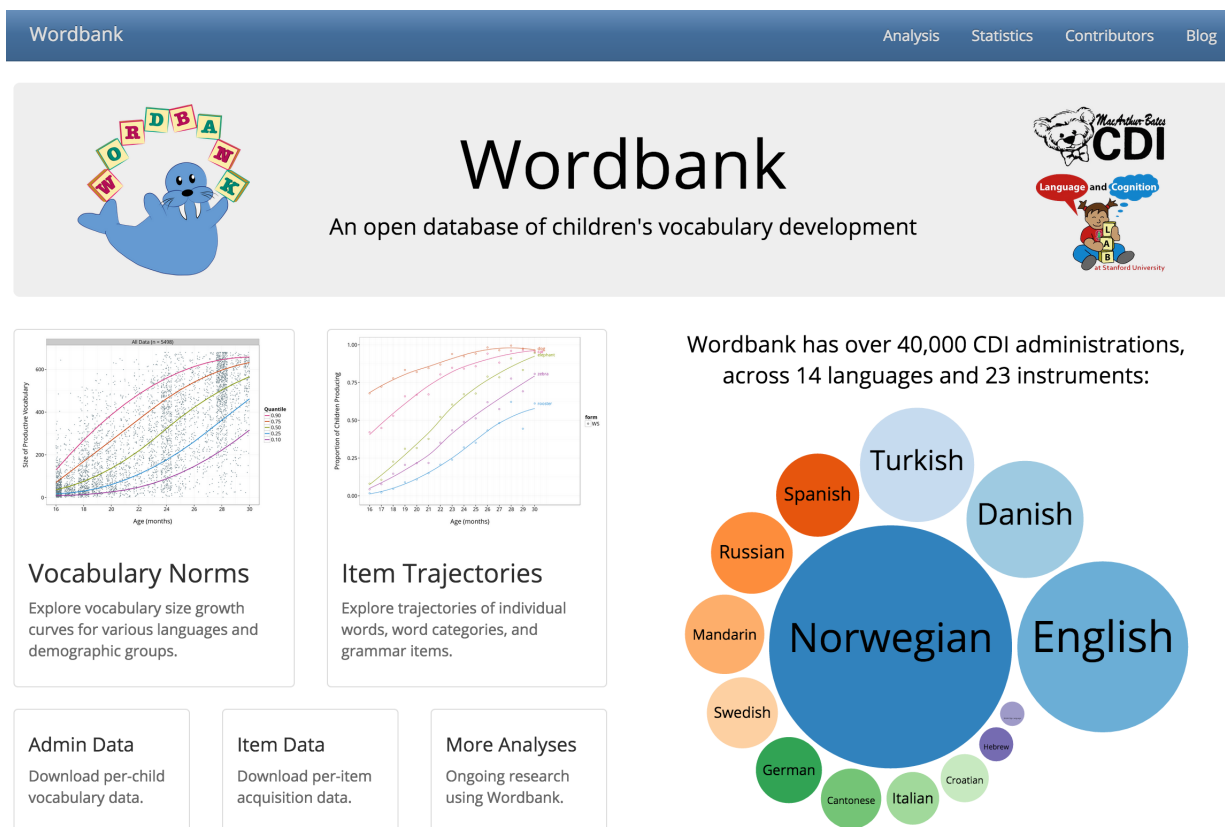
conclusion that a particular child knows a particular word.

Despite these limitations, the CDI instruments are still an important tool when used appropriately. The instruments were designed to minimize bias by targeting current behaviors and asking parents about highly salient features of their child’s abilities. They yield reliable and valid estimates of total vocabulary size, with dozens of studies demonstrating concurrent and predictive relations to naturalistic and observational measures in both typically-developing and at-risk populations (e.g., Dale & Fenson, 1996; Thal, Jackson-Maldonado, & Acosta, 2000; Marchman & Martínez-Sussmann, 2002). In addition, a variety of recent work has shown that individual item-level responses can yield exciting new insights, for example about the growth patterns of semantic networks (Hills et al., 2009; Hills, Maouene, Riordan, & Smith, 2010). Such analyses have the potential to be even more powerful when applied to larger samples and across languages.

## **Wordbank**

To take advantage of the opportunity posed by the broad use of CDI forms in the child language community, we have constructed Wordbank, an open repository for CDI data that allows for interactive analysis and visualization. The current main page of the site is shown in Figure 1. In this section, we begin by describing technical details of the site’s database architecture. We then describe the two primary analysis tools that form the heart of the site’s interactive functionality. We end by discussing the extensibility of the Wordbank framework, highlighting opportunities for contributing data and for building new analyses.

Our inspiration for Wordbank comes from two successful projects for sharing data on children’s language acquisition. The first is the Child Language Data Exchange System (CHILDES; MacWhinney, 2000). A database of transcripts of children’s speech and speech to children, CHILDES has grown into a robust and important tool for the



*Figure 1.* Screenshot of the current Wordbank main page. Visitors can navigate from this page to the interactive reports, as well as to a statistics page that shows the database composition, a contributors page that shows citation information, and a blog that highlights recent updates.

community, with many contributors and affiliated projects. The second is the Cross-Linguistic Lexical Norms site (CLEX; [www.cdi-clex.org/](http://www.cdi-clex.org/); Jørgensen et al., 2010), which is closer in content to Wordbank, and effectively our precursor. CLEX archives normative data from a range of CDI adaptations across languages, allowing browsing of acquisition trajectories for individual items or age groups.

Wordbank builds on CLEX, offering the same functionality but allowing flexible and interactive visualization and analysis, as well as direct database access and data download.



Unique Identifier		Demographic Information						Item-by-Child Data		
A	B	C	D	E	F	G	H	AX	AY	BA
	ender	cdi age	source	birth	month	ethnic	edlev	aabaap	beebah	choochou
1100113		8.00	0.00	1	16	4	2.00	0	0	0
4	1100212	8.00	0.00	1	16	3	2.00	0	0	0
5	1100233	8.00	0.00	2	16	1	2.00	0	0	0
9	1207985	8.00	0.00	2	14	4	2.00	0	0	0
10	1208031	8.00	0.00	3	16	4	2.00	0	1	0

Figure 2. Example data from the CDI norming sample (Fenson et al., 2007). Each row has a unique child identifier, demographics, and word-by-word checklist data.

In addition, Wordbank’s goal is to extend beyond the norming data provided by the developers of individual CDIs by dynamically incorporating data from many different researchers and projects of varying sizes and scopes. While the resulting datasets in Wordbank are likely more heterogeneous, they nevertheless have the potential to be considerably larger and more representative than the individual norming datasets. Wordbank provides tools that enable more powerful, flexible and nuanced analyses of general trends and comparisons across sub-populations in a variety of different languages.

### Database Architecture

Why use a database to store vocabulary data? Consider the standard format of raw CDI data. Figure 2 shows a small slice of the original CDI norming data (Fenson et al., 1994, 2007). Each row is a child, each column gives a variable—either a demographic variable or the result of a particular word being administered to a particular child. Although this format is useful for homogeneous administrations of a single instrument, it cannot accommodate multiple instruments, multiple languages, or datasets with different sources or kinds of demographic information. Consolidating data across different instruments is very difficult in this format, and tracking data on children with multiple longitudinal administrations of a single instrument must also be done in an ad-hoc manner. The move to a database format allows far more flexible and programmatic

handling of heterogeneous data structures from different sources.

A relational database such as Wordbank is at its heart a series of tables linked by unique identifiers. There are two primary groups of tables in Wordbank. The common tables store data that is shared between CDI instruments, including information about children, administrations (individual instances of a form being filled out for a child), and items (words and other questions on a form). The instrument tables store response data for particular CDI instruments.

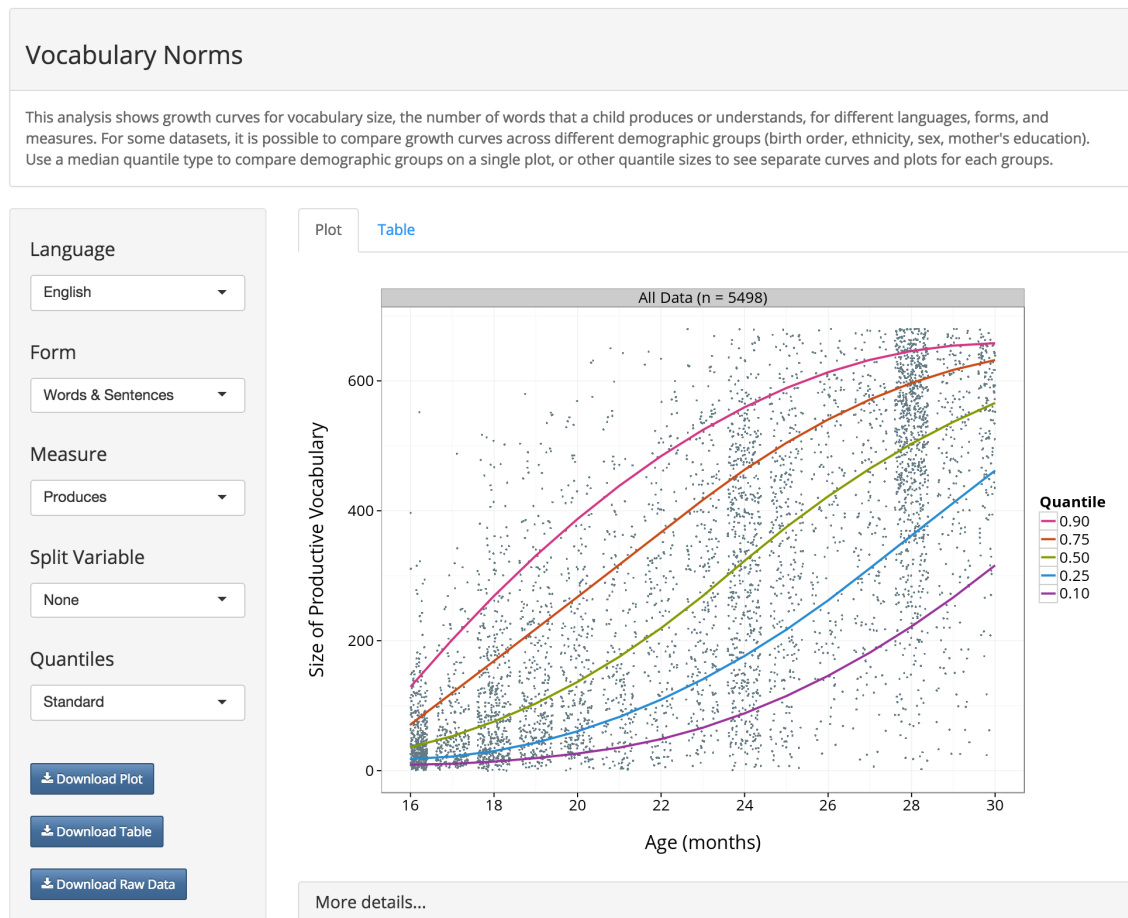
Wordbank is designed so that it can accommodate data from a wide variety of instruments, both within and across languages. Different character sets are supported using unicode. In keeping with the general philosophy of the CDI developers that each new version should be adapted to the unique structure of the new language (up to and including adding new components), our approach to cross-linguistic data is to provide standardized analyses within language, without assuming translation equivalence. At time of writing, the site includes data from more than 44,000 administrations of the CDI across 13 different languages and 23 different instruments.

*Technical details.* Wordbank is constructed using free, open-source tools. The database is a standard MySQL database, managed using Python and Django. Analysis apps are constructed using the **Shiny** package for R, an open-source statistical programming language. The code is hosted in a GitHub repository ([github.com/langcog/wordbank](https://github.com/langcog/wordbank)) where interested users can browse, leave comments, and contribute modifications.

All data uploaded to Wordbank are open and freely available for download, both through the site itself and through the GitHub repository. The site includes only de-identified data that cannot be linked to the parents and children who provided it. Because of these features, the Stanford Institutional Review Board has determined that the Wordbank project does not constitute human subjects research.

### Interactive Analysis Tools

The primary method for users to interact with the Wordbank is through interactive analysis tools that are hosted on the website. These tools allow for fast and flexible exploration of the dataset, the results of which can be exported in tabular and graphical formats for further analysis and presentation.

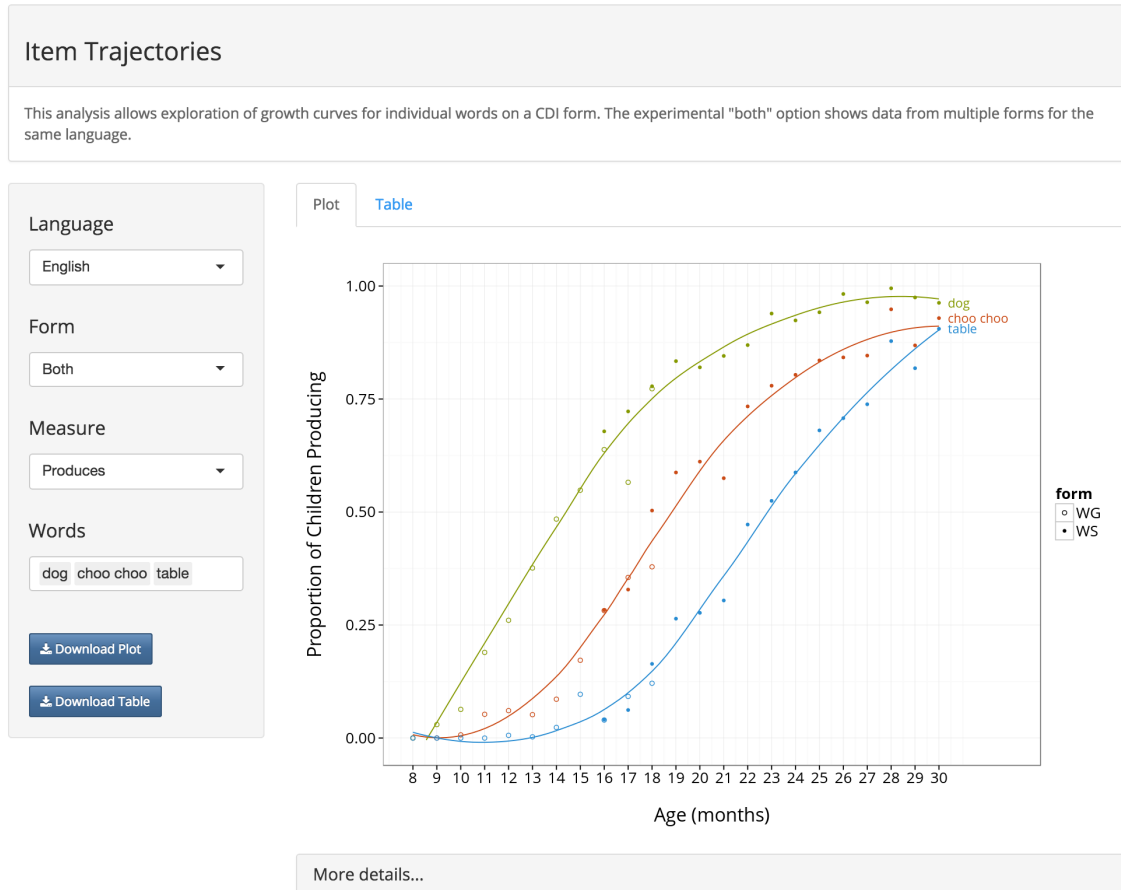


*Figure 3.* A screenshot of the Vocabulary Norms analysis tool, showing 10th, 25th, 50th, 75th, and 90th percentiles (default) for English production scores. Dots show individual administrations, jittered slightly to avoid overplotting. Curves show polynomial spline fits (see text for more details).

*Vocabulary norms.* One of the primary purposes of the CDI form is to provide percentile ranks for vocabulary growth across ages, both for visualizing the variability of early vocabulary growth and for examining differences in these growth patterns due to both individual differences and demographic variables. Accordingly, Wordbank provides a Vocabulary Norms analysis, pictured in Figure 3. The inset plot shows all administrations of a particular CDI instrument within the instrument’s valid age range. Dots show individual children, with age binned by month and jittered to avoid overplotting. Lines on the plot indicate estimates of percentiles, fit using quantile regression with monotonic polynomial splines as the base function (using the `gcrq` package; Muggeo, Sciandra, Tomasello, & Calvo, 2013). An important feature of the norms app is that it can be split by any demographic field in the data, so that comparisons on variables like sex, birth order, or maternal education can be conducted across languages.

The original and updated norming studies (Fenson et al., 1994, 2007) gathered data from a diverse (though not nationally-representative) sample and used these data to construct normative curves from which percentile ranks could be derived. In contrast to these studies, Wordbank is not explicitly designed to provide stable, clinically-relevant norms. Wordbank’s sample is heterogeneous and continually growing, and its analyses are subject to revision and update. Thus, Wordbank does not currently generate percentile ranks for input data, and we do not recommend that Wordbank-generated norming values be used for research or clinical purposes in which the goal is to evaluate children’s performance in reference to an established normative standard. For these types of applications, users should refer to the published norms in the appropriate language.

*Item trajectories.* A second function of the CDI form is to provide aggregate data on the proportion of children at a particular age who know a specific word (Dale & Fenson, 1996; Jørgensen et al., 2010). Such analyses can be extremely helpful for the design and evaluation of materials for young children, including experimental stimuli. Accordingly,



*Figure 4.* A screenshot of the Item Trajectories analysis tool, showing a visualization of the developmental trajectory of production for three words (“dog,” “choo choo,” and “table”) across both Words & Gestures and Words & Sentences forms.

the second major interactive visualization in Wordbank is the Item Trajectories analysis tool.

This tool allows exploration of growth curves for individual words on a CDI form. Users can select a language and instrument (and choose production or comprehension where available), and then select or input a list of words whose trajectories are plotted (Figure 4). The “both” measure option shows data from multiple forms for the same language, with different markers for each item. In general, our exploration suggests that

there are only small differences across different instruments for the same item and age. Lines on the plot show a local polynomial regression smoothing line (`loess` in R).

*Other features: Static reports and tabular data download.* In addition to the interactive analysis tools described above, Wordbank also includes a number of non-interactive but continuously updated reports on features like vocabulary composition across languages, links between grammar and the lexicon (Braginsky, Yurovsky, Marchman, & Frank, 2015), and gender differences in vocabulary growth (see below). On the Analysis page ([wordbank.stanford.edu/reports](http://wordbank.stanford.edu/reports)), we provide a gallery of both interactive and non-interactive analyses.

Wordbank also allows raw tabular data to be browsed and downloaded for subsequent analysis in all popular statistical packages. Using the same basic interface as the Vocabulary Norms and Item Trajectory tools, users can browse raw data aggregated across children (similar to the Vocabulary Norms tool), across items (similar to the Items Trajectory tool), or even view the raw subject-by-item data. All data in these “standard” reports can be downloaded in CSV format.

### *Extensibility*

Extensibility is one of the major strengths of Wordbank. Although programming knowledge is not necessary for interacting with Wordbank, interested researchers with programming skills can contribute to the development effort by adding new analyses. Each Wordbank analysis app is constructed as a standalone script or set of scripts in the R language. Constructing an interactive analysis requires specifying a visualization and some interactive functionality using **Shiny**. Non-interactive analyses can be constructed as R **Markdown** documents that execute scripts using the Wordbank database. Both of these have the virtue of rerunning on the newest version of the database whenever they are opened, so they do not go out of date as new data are added.

In addition, we encourage contributions of individual datasets. Wordbank currently imports data from Excel and CSV formats via automated import scripts. Individuals or labs interested in contributing should consult with the authors for advice about data formatting and upload.

### **wordbankr: an R package for accessing Wordbank**

Although the analysis tools described above suffice for many needs, researchers interested in detailed quantitative or cross-linguistic analyses may wish to connect directly to the Wordbank database and manipulate the data directly. To facilitate this functionality, we provide the **wordbankr** package for the popular R programming language. This package abstracts away the details of connecting to the database. Users can take advantage of the SQL tools developed in the popular **dplyr** package (Wickham & Francois, 2014), which make manipulating large datasets quick and easy. We describe the commands that the package provides and then give a worked example of using the package for a simple analysis.

#### *Package details*

The **wordbankr** package provides a function for connecting to the Wordbank database, **connect\_to\_wordbank**. Users can connect in **remote** mode, which loads data from the Wordbank server, or in **local** mode if they have a copy of the database set up on their local machine. This connection can then be used to load any of the common tables by name with **get\_common\_table** and any of the instrument tables by language and form with **get\_instrument\_table**. The package also provides functions that make loading components of the data easier by joining together relevant tables to give different views of the data: by-administration (**get\_administration\_data**) with demographics and vocabulary sizes; by-item (**get\_item\_data**) with item types, categories and other information; or administration-by-item (**get\_instrument\_data**) with raw response values.

For more detailed documentation, see the package repository (<http://github.com/langcog/wordbankr>).

*A worked example: Sex differences across languages*

We demonstrate the analytic potential of direct manipulation of the Wordbank database using **wordbankR**. Our example is an analysis examining sex differences. Sex differences in productive language are commonly found in individual studies (e.g., Fenson et al., 1994; Huttenlocher, Haight, Bryk, Seltzer, & Lyons, 1991; see Wallentin, 2009 for review), and one large-scale previous study found differences in productive vocabulary in 10 languages (Eriksson et al., 2012). In the following analysis, we show how this general analysis can quickly be replicated using Wordbank.

To perform the analysis, we first begin by using **wordbankr** to load the data from Wordbank and connect to the tables we need:

```
admins <- get_administration_data()
items <- get_item_data()
```

We next use a series of **dplyr** calls to compute the number of words in each language, select the appropriate subset of the data, and calculate the proportion of words produced for this data subset:

```
num_words <- items %>%
  filter(form == "WS", type == "word") %>%
  group_by(language) %>%
  summarise(n = n())

vocab_admins <- admins %>%
  select(data_id, language, form, age, sex, production)

vocab_data <- vocab_admins %>%
```



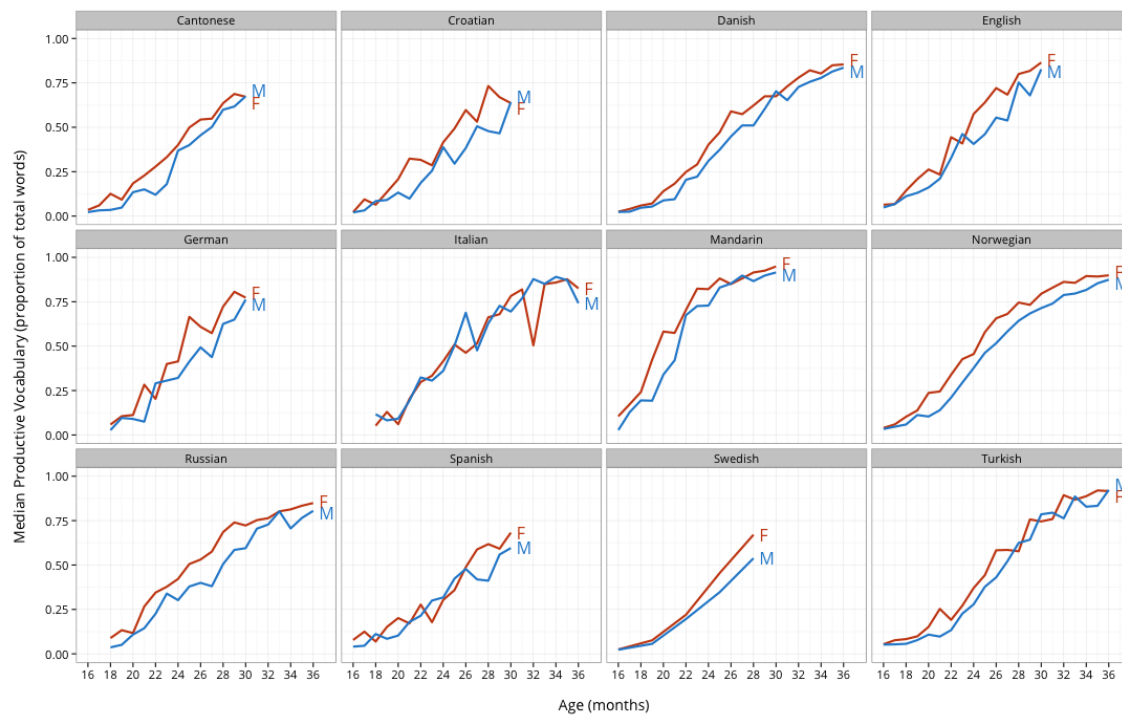
```

group_by(language, sex, age) %>%
left_join(num_words) %>%
mutate(production = as.numeric(production) / n) %>%
summarise(median = median(production))

```

We then plot the `vocab_data` data frame using the `ggplot2` library (Wickham, 2009).

Full code for the analysis as a whole (including the plot) is given in the Using Wordbank tutorial, available on at [wordbank.stanford.edu/tutorial](http://wordbank.stanford.edu/tutorial).



*Figure 5.* Median productive vocabulary as a proportion of total words on an instrument, plotted by age in months. Red and blue lines show females and males, respectively.

The results of this analysis are shown in Figure 5. As expected, we replicate the gender differences found in previous work: Female showed a small but highly reliable advantage in early production. This effect is highly consistent and clearly visible in 11 out

of 12 languages, with Italian being the only exception. While the Wordbank data do not allow us to speculate about the origins of this difference, they certainly allow us to formulate hypotheses with substantially more clarity than previous analyses. Such visualizations also highlight differences in the size and composition of the database.

## Conclusions

In this paper, we presented Wordbank, an open repository for parent-report vocabulary data using the MacArthur-Bates CDI. The interactive analysis tools available on the Wordbank site allow interested researchers to explore a wide variety of phenomena in vocabulary development quickly and easily, exporting data and downloading presentation-quality graphics that document their analysis. In addition, users can contribute new analyses and data to the site and connect to it directly using an R package for data loading. These functions all enable greater sharing and reuse of existing data on children’s vocabulary. We hope that the resulting tools will enable new discoveries in the future.

## References

- Bates, E. (1976). *Language and context: The acquisition of pragmatics* (Vol. 13). New York, NY: Academic Press.
- Bates, E., & Goodman, J. (1999). On the emergence of grammar from the lexicon. In B. Macwhinney (Ed.), *The emergence of language*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P., Reznick, J. S., . . . Hartung, J. (1994). Developmental and stylistic variation in the composition of early vocabulary. *Journal of Child Language*, 21, 85–123.
- Bloom, P. (2002). *How children learn the meanings of words*. Cambridge, MA: MIT Press.

- Bornstein, M. H., & Haynes, O. M. (1998). Vocabulary competence in early childhood: Measurement, latent construct, and predictive validity. *Child Development*, 69, 654–671.
- Braginsky, M., Yurovsky, D., Marchman, V. A., & Frank, M. C. (2015). Developmental changes in the relationship between grammar and the lexicon. In *Proceedings of the 37th Annual Meeting of the Cognitive Science Society*.
- Brown, R. (1973). *A first language: The early stages*. Cambridge, MA: Harvard University Press.
- Dale, P. S., & Fenson, L. (1996). Lexical development norms for young children. *Behavior Research Methods, Instruments, & Computers*, 28, 125–127.
- Dunn, L. M., & Dunn, L. M. (2007). *Peabody picture vocabulary test* (4th Edition ed.). Parsippany, NJ: AGS Publishing/Pearson Assessments.
- Eriksson, M., Marschik, P. B., Tulviste, T., Almgren, M., Pérez Pereira, M., Wehberg, S., ... Gallego, C. (2012). Differences between girls and boys in emerging language skills: Evidence from 10 language communities. *British Journal of Developmental Psychology*, 30, 326–343.
- Feldman, H. M., Dale, P. S., Campbell, T. F., Colborn, D. K., Kurs-Lasky, M., Rockette, H. E., & Paradise, J. L. (2005). Concurrent and predictive validity of parent reports of child language at ages 2 and 3 years. *Child Development*, 76, 856–868.
- Feldman, H. M., Dollaghan, C. A., Campbell, T. F., Kurs-Lasky, M., Janosky, J. E., & Paradise, J. L. (2000). Measurement properties of the macarthur communicative development inventories at ages one and two years. *Child Development*, 71, 310–322.
- Fenson, L., Bates, E., Dale, P., Goodman, J., Reznick, J. S., & Thal, D. (2000). Reply: Measuring variability in early child language: Don't shoot the messenger. *Child Development*, 71, 323–328.

- Fenson, L., Dale, P., Reznick, J., Bates, E., Thal, D., Pethick, S., . . . Stiles, J. (1994). Variability in early communicative development. *Monographs of the Society for Research in Child Development*, 59.
- Fenson, L., Marchman, V. A., Thal, D., Dale, P., Reznick, J. S., & Bates, E. (2007). *MacArthur-Bates Communicative Development Inventories: User's Guide and Technical Manual* (2nd ed.). Baltimore, MD: Brookes Publishing Company.
- Hart, B., & Risley, T. (1995). *Meaningful differences in the everyday experience of young american children*. Baltimore, MD: Brookes Publishing Company.
- Hills, T. T., Maouene, J., Riordan, B., & Smith, L. B. (2010). The associative structure of language: Contextual diversity in early word learning. *Journal of Memory and Language*, 63, 259–273.
- Hills, T. T., Maouene, M., Maouene, J., Sheya, A., & Smith, L. (2009). Longitudinal analysis of early semantic networks preferential attachment or preferential acquisition? *Psychological Science*, 20, 729–739.
- Huttenlocher, J., Haight, W., Bryk, A., Seltzer, M., & Lyons, T. (1991). Early vocabulary growth: Relation to language input and gender. *Developmental Psychology*, 27, 236–248.
- Jørgensen, R. N., Dale, P. S., Bleses, D., & Fenson, L. (2010). CLEX: A cross-linguistic lexical norms database. *Journal of Child Language*, 37, 419–428.
- Kristoffersen, K. E., Simonsen, H. G., Bleses, D., Wehberg, S., Jørgensen, R. N., Eiesland, E. A., & Henriksen, L. Y. (2013). The use of the internet in collecting cdi data—an example from norway. *Journal of Child Language*, 40, 567–585.
- MacWhinney, B. (2000). *The CHILDES Project: Tools for Analyzing Talk. Third Edition*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Marchman, V. A., & Fernald, A. (2008). Speed of word recognition and vocabulary knowledge in infancy predict cognitive and language outcomes in later childhood.

*Developmental Science*, 11, F9–F16.

- Marchman, V. A., & Martínez-Sussmann, C. (2002). Concurrent validity of caregiver/parent report measures of language for children who are learning both english and spanish. *Journal of Speech, Language, and Hearing Research*, 45, 983–997.
- Muggeo, V. M., Sciandra, M., Tomasello, A., & Calvo, S. (2013). Estimating growth charts via nonparametric quantile regression: a practical framework with application in ecology. *Environmental and Ecological Statistics*, 20, 519–531.
- Nelson, K. (1973). Structure and strategy in learning to talk. *Monographs of the Society for Research in Child Development*, 1–135.
- Rescorla, L. (1989). The language development survey: a screening tool for delayed language in toddlers. *Journal of Speech and Hearing Disorders*, 54, 587–599.
- Tardif, T., Fletcher, P., Liang, W., Zhang, Z., Kaciroti, N., & Marchman, V. A. (2008). Baby’s first 10 words. *Developmental Psychology*, 44, 929.
- Thal, D., Jackson-Maldonado, D., & Acosta, D. (2000). Validity of a parent-report measure of vocabulary and grammar for spanish-speaking toddlers. *Journal of Speech, Language, and Hearing Research*, 43, 1087–1100.
- Tomasello, M., & Mervis, C. B. (1994). The instrument is great, but measuring comprehension is still a problem. *Monographs of the Society for Research in Child Development*, 59, 174–179.
- Wallentin, M. (2009). Putative sex differences in verbal abilities and language cortex: A critical review. *Brain and Language*, 108, 175–183.
- Wickham, H. (2009). *ggplot2: elegant graphics for data analysis*. Springer Science & Business Media.
- Wickham, H., & Francois, R. (2014). dplyr: A grammar of data manipulation. *R package version 0.3.0.2*.