1	childes-db: a flexible and reproducible interface to the Child Language Data Exchange
2	System
3	Alessandro Sanchez ^{*1} , Stephan C. Meylan ^{*2} , Mika Braginsky ³ , Kyle E. MacDonald ¹ , Daniel Yurovsky ⁴ , & Michael C. Frank ¹
5	¹ Stanford University
6	² University of California, Berkeley
7	3 MIT
8	⁴ University of Chicago

Author Note
 Thanks to Brian MacWhinney for advice and guidance, and to Melissa Kline for her
 work on ClanToR, which formed a starting point for our work. This work is supported by a
 Jacobs Advanced Research Fellowship to MCF.
 Correspondence concerning this article should be addressed to Alessandro Sanchez*,

¹⁴ Department of Psychology, 450 Serra Mall, Stanford, CA 94305. E-mail:

15 sanchez7@stanford.edu

16

Abstract

The Child Language Data Exchange System (CHILDES) has played a critical role in 17 research on child language development, particularly in characterizing the early language 18 learning environment. Access to these data can be both complex for novices and difficult to 19 automate for advanced users, however. To address these issues, we introduce childes-db, a 20 database-formatted mirror of CHILDES that improves data accessibility and usability by 21 offering novel interfaces, including browsable web applications and an R application 22 programming interface (API). Along with versioned infrastructure that facilitates 23 reproducibility of past analyses, these interfaces lower barriers to analyzing naturalistic 24 parent-child language, allowing for a wider range of researchers in language and cognitive 25 development to easily leverage CHILDES in their work. 26

Keywords: child language; corpus linguistics; reproducibility; R packages; research
 software

²⁹ Word count: 2925

childes-db: a flexible and reproducible interface to the Child Language Data Exchange
 System

32

Introduction

What are the representations that children learn about language, and how do they 33 emerge from the interaction of learning mechanisms and environmental input? Developing 34 facility with language requires learning a great many interlocking components – meaningful 35 distinctions between sounds (phonology), names of particular objects and actions (word 36 learning), meaningful sub-word structure (morphology), rules for how to organize words 37 together (syntax), and context-dependent and context-independent aspects of meaning 38 (semantics and pragmatics). Key to learning all of these systems is the contribution of the 39 child's input – exposure to linguistic and non-linguistic data – in the early environment. 40 While in-lab experiments can shed light on linguistic knowledge and some of the implicated 41 learning mechanisms, characterizing this early environment requires additional research 42 methods and resources. 43

One of the key methods that has emerged to address this gap is the collection and 44 annotation of speech to and by children, often in the context of the home. Starting with 45 Roger Brown's (1973) work on Adam, Eve, and Sarah, audio recordings – and more recently 46 video recordings – have been augmented with rich, searchable annotations to allow 47 researchers to address a number of questions regarding the language learning environment. 48 Focusing on language learning in naturalistic contexts also reveals that children have, in 49 many cases, productive and receptive abilities exceeding those demonstrated in experimental 50 contexts. Often, children's most revealing and sophisticated uses of language emerge in the 51 course of naturalistic play. 52

⁵³ While corpora of early language acquisition are extremely useful, creating them ⁵⁴ requires significant resources. Collecting and transcribing audio and video is costly and ⁵⁵ extremely time consuming – even orthographic transcription (i.e., transcriptions with ⁵⁶ minimal phonetic detail) can take ten times the duration of the original recording

CHILDES-DB: AN INTERFACE TO CHILDES

(MacWhinney, 2000). Automated, machine learning-based methods like automatic speech recognition (ASR) have provided only modest gains in efficiency. Such systems are limited both by the less-than-ideal acoustic properties of home recordings, and also by the poor fit of language models built on adult-directed, adult-produced language samples to child-directed and child-produced speech. Thus, researchers' desires for data in analyses of child language corpora can very quickly outstrip their resources.

Established in 1984 to address this issue, the Child Language Data Exchange System 63 (CHILDES) aims to make transcripts and recordings relevant to the study of child language 64 acquisition available to researchers as free, public datasets (MacWhinney, 2000, 2014; 65 MacWhinney & Snow, 1985). CHILDES now archives tens of thousands of transcripts and 66 associated media across 20+ languages, making it a critical resource for characterizing both 67 children's early productive language use and their language environment. As the first major 68 effort to consolidate and share transcripts of child language, CHILDES has been a pioneer in 69 the move to curate and disseminate large-scale behavioral datasets publicly. 70

Since its inception, a tremendous body of research has made use of CHILDES data. 71 Individual studies are too numerous to list, but classics include studies of morphological 72 over-regularization (Marcus et al., 1992), distributional learning (Redington, Chater, & 73 Finch, 1998), word segmentation (Goldwater, Griffiths, & Johnson, 2009), the role of 74 frequency in word learning (Goodman, Dale, & Li, 2008), and many others. Some studies 75 analyze individual examples in depth (e.g., Snyder, 2007), others track multiple 76 child-caregiver dyads (e.g., Meylan, Frank, Roy, & Levy, 2017), and still others use the 77 aggregate properties of all child or caregiver speech pooled across corpora (Montag, Jones, & 78 Smith, 2015; e.g., Redington et al., 1998). 79

Nonetheless, there are some outstanding challenges working with CHILDES, both for
 students and for advanced users. The CHILDES ecosystem uses a specialized file format
 (CHAT), which is stored as plain text but includes structured annotations grouped into tiers
 stored on separate lines. These tiers allow information about utterances to be stored with

accompanying information such as the phonological, morphological, or syntactic structure of
the utterance. These files are usually analyzed using a command-line program (CLAN) that
allows users to count word frequencies, compute statistics (e.g., mean length of utterance, or
MLU), and execute complex searches against the data. While this system is flexible and
powerful, mastering the CHAT codes and especially the CLAN tool with its many functions
and flags can be daunting. These technical barriers decrease the ease of exploration by a
novice researcher or in a classroom exercise.

On the opposite end of the spectrum, for data-oriented researchers who are interested in doing large-scale analyses of CHILDES, the current tools are also not ideal. CLAN software is an excellent tool for interactive exploration, but – as a free-standing application – it can be tricky to build into a processing pipeline written in Python or R. Thus, researchers who would like to ingest the entire corpus (or some large subset) into a computational analysis typically write their own parsers of the CHAT format to extract the subset of the data they would like to use (Meylan et al., 2017; e.g., Redington et al., 1998; Yang, 2013).

The practice of writing custom parsers is problematic for a number of reasons. First, 98 effort is wasted in implementing the same features again and again. Second, this process can 99 introduce errors and inconsistencies in data handling due to difficulties dealing with the 100 many special cases in the CHAT standard. Third, these parsing scripts are rarely shared – 101 and when when they are, they typically break with subsequent revisions to the dataset – 102 leading to much greater difficulty in reproducing the exact numerical results from previous 103 published research that used CHILDES (see e.g., Meylan et al., 2017 for an example). 104 Fourth, the CHILDES corpus itself is a moving target: computational work using the entire 105 corpus at one time point may include a different set of data than subsequent work due as 106 corpora are added and revised. Currently, there is no simple way for researchers to document 107 exactly which version of the corpus has been used, short of creating a full mirror of the data. 108 These factors together lead to a lack of *computational reproducibility*, a major problem that 109 keeps researchers from verifying or building on published research (Donoho, 2010; Stodden et 110

¹¹¹ al., 2016).

In the current manuscript, we describe a system for extending the functionality of 112 CHILDES to address these issues. Our system, childes-db, is a database-formatted mirror 113 of CHILDES that allows access through an application programming interface (API). This 114 infrastructure allows the creation of web applications for browsing and easily visualizing the 115 data, facilitating classroom use of the dataset. Further, the database can be accessed 116 programmatically by advanced researchers, obviating the need to write one-off parsers of the 117 CHAT format. The database is versioned for access to previous releases, allowing 118 computational reproducibility of particular analyses. 119

We begin by describing the architecture of childes-db and the web applications that we provide. Next, we describe the childesr API, which provides a set of R functions for programmatic access to the data while abstracting away many of the technical details. We conclude by presenting several worked examples of specific uses of the system – both web apps and the R API – for research and teaching.

125

Design and technical approach

As described above, CHILDES is most often approached as a set of distinct CHAT 126 files, which are then parsed by users, often using CLAN. In contrast to this parsing approach, 127 which entails the sequential processing of strings, childes-db treats CHILDES as a set of 128 linked tables, with records corresponding to intuitive abstractions such as words, utterances, 129 and transcripts (see Kline, 2012 for an earlier example of deriving tabular representations of 130 CHILDES). Users of data analysis languages like R or Julia, libraries like Pandas, or those 131 familiar with Structured Query Language (SQL) will be familiar with operations on tables 132 such as filtering (subsetting), sorting, aggregation (grouping), and joins (merges). These 133 operations obviate the need for users to consider the specifics of the CHAT representation – 134 instead they simply request the entities they need for their research and allow the API to 135 take care of the formatting details. We begin by orienting readers to the design of the system 136



Figure 1. Database schema for 'childes-db'. Tokens are linked to superordinate groupings of utterances, transcripts, corpora, and collections (red arrows). All tokens and utterances are additionally associated with a participant (blue arrows).

via a top-level description and motivation for the design of the database schema, then
provide details on the database's current technical implementation and the versioning
scheme. Users primarily interested in accessing the database can skip these details and focus
on access through the childesr API and the web apps.

141 Database format

At its core, childes-db is a database consisting of a set of linked tabular data stores 142 where records correspond to linguistic entities like words, utterances, and sampling units like 143 transcriptions and corpora. The smallest unit of abstraction tracked by the database is a 144 token, treated here as the standard (or citation) orthographic form of a word. Using the 145 standardized written form of the word facilitates the computation of lexical frequency 146 statistics for comparison or aggregation across children or time periods. Deviations from the 147 citation form – which are particularly common in the course of language development and 148 often of interest to researchers – are kept as a separate (possibly null) field associated with 149 each token. 150

Many of the other tables in the database are hierarchical collections built out of tokens *utterance, transcript, corpus,* and *collection* – that store attributes appropriate for each level of description. Every entity includes attributes that link it to all higher-order collections, *e.g.,* an utterance lists the transcript, corpus, and collection to which it belongs. An

utterance contains one or more words and includes fields such as the utterance type such as 155 *declarative* or *interrogative*, total number of tokens, and the total number of morphemes if 156 the morphological structure is available in the original CHAT file. A transcript consists of 157 one or more utterances and includes the date collected, the name of the target child, and age 158 in days if defined, and the filename from CHILDES. A *corpus* consists of one or more 159 transcripts, corresponding to well-known collections like the Brown (Brown, 1973) or 160 Providence (Demuth, Culbertson, & Alter, 2006) corpora. Finally, a *collection* is a 161 superordinate collection of corpora generally corresponding to a geographic region, following 162 the convention in CHILDES. Because every record can be linked to a top-level collection 163 (generally corresponding to a language), each table includes data from all languages 164 represented in CHILDES. 165

Participants – generally children and caregivers – are represented separately from the 166 token hierarchy because it is common for the same children to appear in multiple transcripts. 167 A participant identifier is associated with every word and utterance, including a name, role, 168 3-letter CHILDES identifier (CHI = child, MOT = mother, FAT = father, etc.), and the 169 range of ages (or age of corresponding child) for which they are observed. For non-child 170 participants (caregivers and others), the record additionally contains an identifier for the 171 corresponding target child, such that data corresponding to children and their caregivers can 172 be easily associated. 173

174 Technical implementation

childes-db is stored as a MySQL database, an industry-standard, open-source
relational database server that can be accessed directly from a wide range of programming
languages. The childes-db project provides hosted, read-only databases for direct access and
for childesr (described below) as well as compressed .sql exports for local installation.
While the former is appropriate for most users, local installation can provide performance
gains by allowing a user to access the database on their machine or on their local network, as

¹⁸¹ well as allowing users to store derived information in the same database.

In order to import the CHILDES corpora into the MySQL schema described above, it must first be accurately parsed and subsequently vetted to ensure its integrity. We parse the XML (eXtensible Markup Language) release of CHILDES hosted by childes.talkbank.org using the NLTK library in Python (Bird & Loper, 2004). Logic implemented in Python converts the linear, multi-tier parse into a tabular format appropriate for childes-db. This logic includes decisions that we review below regarding what information sources are captured in the current release of the database and which are left for future development.

The data imported into childes-db is subject to data integrity checks to ensure that 189 our import of the corpora is accurate and preferable over ad-hoc parsers developed by many 190 individual researchers. In order to evaluate our success in replicating CLAN parses, we 191 compared unigram counts in our database with those outputted by CLAN, the 192 command-line tool built specifically for analysis of transcripts coded in CHAT. We used the 193 CLAN commands FREQ and MLU to compare total token counts and mean lengths of 194 utterance for every speaker in every transcript and compared these these values to our own 195 using the Pearson correlation coefficient. The results of the comparison were .99 and .98 for 196 the unigram count and MLU data, respectively, indicating reliable parsing. 197

Versioning. The content of CHILDES changes as additional corpora are added or 198 transcriptions are updated; as of time of writing, these changes are not systematically 199 tracked. To facilitate reproducibility of past analyses, we introduce a simple versioning 200 system by adding a new complete parse of the current state of CHILDES every six months 201 or as warranted by changes in CHILDES. By default, users interact with the most recent 202 version of the database available. To support reproduction of results with previous versions 203 of the database, we continue to host recent versions (up to the last three years / six versions) 204 through our childesr API so that researchers can run analyses against specific historical 205 versions of the database. For versions more than three years old, we host compressed .sql 206 files that users may download and serve using a local installation of MySQL server. 207

CHILDES-DB: AN INTERFACE TO CHILDES

Current Annotation Coverage. The current implementation of childes-db
emphasizes the computation of lexical statistics, and consequently focuses on reproducing
the words, utterances, and speaker information in CHILDES transcripts. For this reason, we
do not preserve all of the information available in CHILDES, such as:

• Sparsely annotated tiers, e.g. phonology (%pho) and situation (%sit)

• Media links

• Tone direction and stress

• Filled pauses

• Reformulations, word revision, and phrase revision, e.g. <what did you>[//] how can you see it ?

• paralinguistic material, e.g. [=! cries]

We will prioritize the addition of these information sources and others in response to community feedback.

221

Interfaces for Accessing childes-db

We first discuss the childes-db web apps and then introduce the childesr R package.

223 Interactive Web Apps

The ability to easily browse and explore the CHILDES corpora is a cornerstone of the childes-db project. To this end we have created powerful yet easy-to-use interactive web applications that enable users to visualize various dimensions of the CHILDES corpus: frequency counts, mean lengths of utterance, type-token ratios, and more. All of this is doable without the requirement of understanding command-line tools or any kind of programming knowledge as had been the case with CLAN.¹

¹ The LuCiD toolkit (Chang, 2017) provides related functionality for a number of common analyses. In contrast to those tools, which focus on filling gaps not covered by CLAN – e.g., the use of n-gram models,

Our web apps are built using Shiny, an R package that enables easy app construction 230 using R. Underneath the hood, each web app is making calls to our childesr API and 231 subsequently plots the data using the popular R plotting package ggplot2. A user's only 232 task is to configure exactly what should be plotted through a series of buttons, sliders, and 233 text boxes. The user may specify what collection, corpus, child, age range, caregiver, etc., 234 should be included in a given analysis. The plot is displayed and updated in real-time, and 235 the underlying data are also available for download alongside the plot. All of these analyses 236 may also be reproduced using the childesr package, but the web apps are intended for the 237 casual user who seeks to easily extract developmental indices quickly and without any 238 technical overhead. 239

Frequency Counts. The lexical statistics of language input to children have long 240 been an object of study in child language acquisition research. Frequency counts of words in 241 particular may provide insight into the cognitive, conceptual, and linguistic experience of a 242 young child (see e.g., Ambridge, Kidd, Rowland, & Theakston, 2015 for review). In this web 243 app, inspired by ChildFreq (Bååth, 2010), we provide users the ability to search for any word 244 spoken by a participant in the CHILDES corpora and track the usage of that word by a 245 child or caregiver over time. Because of the various toggles available to the user that can 246 subset the data, a user may word frequencies curves for a single child in the Brown corpus or 247 all Spanish speaking children, if desired. In addition, users can plot frequency curves 248 belonging to caregivers alongside their child for convenient side-by-side comparisons. A 240 single word or multiple words may be entered into the input box. 250

Derived Measures. The syntactic complexity and lexical diversity of children's speech are similarly critical metrics for acquisition researchers (Miller & Chapman, 1981; Watkins, Kelly, Harbers, & Hollis, 1995). There are a number of well-established measures of children's speech that operationalize complexity and diversity, and have many applications in

incremental sentence generation, and distributional word classification – our web apps focus on covering the same common tasks as CLAN, but making the outputs into browsable visualizations.



Figure 2. Frequency Counts.

speech-language pathology (SLP), where measures outside of the normal range may be
indicative of speech, language, or communication disorders.

Several of the most common of these measures are available in the Derived Measures app, which plots these measures across age for a given subset of data, again specified by collection, corpora, children, and speakers. As with the Frequency Counts app, caregivers' lexical diversity measures can be plotted alongside children's.

We have currently implemented the following measures:

- MLU-w (mean length of utterance in words),
- MLU-m (mean length of utterance in morphemes),
- TTR (type-token ratio, a measure of lexical diversity; Templin, 1957),
- MTLD (measure of textual lexical diversity; Malvern & Richards, 1997),
- HD-D (lexical diversity via the hypergeometric distribution; McCarthy & Jarvis, 2010)

²⁶⁷ As with the Frequency Counts app, a user may subset the data as they choose, compare

measures between caregivers and children, and aggregate across children from different

269 corpora.



Figure 3. Derived Measures.

Population Viewer. Many times a researcher will want to investigate the statistics of corpora (e.g., their size, number of utterances, number of tokens) before choosing a target corpus or set of corpora for a project. This web app is intended to provide a basic overview regarding the scale and temporal extent of various corpora in CHILDES, as well as giving researchers insight into the aggregate characteristics of CHILDES. For example, examining the aggregate statistics reveals that coverage in CHILDES peaks at around 30 months.

276 The childesr Package

Although the interactive analysis tools described above cover some of the most common use cases of CHILDES data, researchers interested in more detailed and flexible analyses will want to interface directly with the data in childes-db. Making use of the R programming



Figure 4. Population Viewer.

language (R Core Team, 2017), we provide the childesr package. R is an open-source, 280 extensible statistical computing environment that is rapidly growing in popularity across 281 fields and is increasing in use in child language research (Norrman & Bylund, 2015; e.g. 282 Song, Shattuck-Hufnagel, & Demuth, 2015). The childesr package abstracts away the 283 details of connecting to and querying the database. Users can take advantage of the tools 284 developed in the popular dplyr package (Wickham, Francois, Henry, & Müller, 2017), which 285 makes manipulating large datasets quick and easy. We describe the commands that the 286 package provides and then give several worked examples of using the package for analyses. 287 The childesr package is easily installed via CRAN, the comprehensive R archive 288 network. To install, simply type: install.packages("childesr"). After installation, users 280 have access to functions that can be used to retrieve tabular data from the database: 290

291

• get_collections() gives the names of available collections of corpora ("Eng-NA",

²⁹² "Spanish", etc.)

• get_corpora() gives the names of available corpora ("Brown", "Clark", etc.)

• get_transcripts() gives information on available transcripts (language, date, target

295	child demographics)
296 •	get_participants() gives information on transcript participants (name, role,
297	demographics)
298 •	get_speaker_statistics() gives summary statistics for each participant in each
299	transcript (number of utterances, number of types, number of tokens, mean length of
300	utterance)
301 •	get_utterances() gives information on each utterance (glosses, stems, parts of
302	speech, utterance type, number of tokens, number of morphemes, speaker information
303	target child information)
304 •	get_types() gives information on each type within each transcript (gloss, count,
305	speaker information, target child information)

get_tokens() gives information on each token (gloss, stem, part of speech, number of
 morphemes, speaker information, target child information)

Each of these functions take arguments that restrict the query to a particular subset of the data (e.g. by collection, by corpus, by speaker role, by target child age, etc.) and returns the output in the form of a table. All functions support the specification of the database version to use. For more detailed documentation, see the package repository

312 (http://github.com/langcog/childesr).

313

Using childes-db: Worked Examples

In this section we give a number of examples of how childes-db can be used in both research and teaching, using both the web apps and the R API. Note that all of these examples use dplyr syntax (Wickham et al., 2017); several accessible introductions to this framework are available online (e.g., Wickham & Grolemund, 2016).

318 Research applications

Color frequency. One common use of CHILDES is to estimate the frequency with 319 which children hear different words. These frequency estimates are used both in the 320 development of theory (e.g., frequent words are learned earlier; Goodman et al., 2008), and 321 in the construction of age-appropriate experimental stimuli. One benefit of the childes-db 322 interface is that it allows for easy analysis of how the frequencies of words change over 323 development. Many of our theories in which children learn the structure of language from its 324 statistical properties implicitly assume that these statistics are *stationary*, i.e. unchanging 325 over development (e.g., Saffran, Aslin, & Newport, 1996). However a number of recent 326 analyses show that the frequencies with which infants encounter both linguistic and visual 327 properties of their environment may change dramatically over development (Fausey, 328 Jayaraman, & Smith, 2016), and these changing distributions may produce similarly 329 dramatic changes in the ease or difficulty with which these regularities can be learned 330 (Elman, 1993). 331

To demonstrate how one might discover such non-stationarity, we take as a case study 332 the frequency with which children hear the color words of English (e.g. "blue", "green"). 333 Color words tend to be learned relatively late by children, potentially in part due to the 334 abstractness of the meanings to which they refer (see Wagner, Dobkins, & Barner, 2013). 335 However, within the set of color words, the frequency with which these words are heard 336 predicts a significant fraction of the variance in their order of acquisition (Yurovsky, Wagner, 337 Barner, & Frank, 2015). But are these frequencies stationary – e.g. do children hear "blue" 338 as often at 12 months as they do at 24 months? We answer this question in two ways – first 339 using the web apps, and then using the childesr package. 340

³⁴¹ Using web apps. To investigate whether the frequency of color words is stationary ³⁴² over development, a user can navigate to the Frequency app, and enter a set of color words ³⁴³ into the Word selector separated by a comma: here "blue, red, green." Because the question ³⁴⁴ of interest is about the frequency of words in the input (rather than produced by children), ³⁴⁵ the Speaker field can be set to reflect this choice. In this example we select "Mother." Because children learn most of their basic color words by the age of 5, the age range 1–5 years is a reasonable choice for Ages to include. The results of these selections are shown in Figure 5. We can also create a hyperlink to store these set of choices so that we can share these results with others (or with ourselves in the future) by clicking on the Share Analysis button in the bottom left corner.

From this figure, it seems likely that children hear "blue" more frequently early in development, but the trajectories of "red" and "green" are less clear. We also do not have a good sense of the errors of these measurements, are limited to just a few colors at a time before the plot becomes too crowded, and cannot combine frequencies across speakers. To perform this analysis in a more compelling and complete way, a user can use the childesr interface.



Figure 5. An example of using the Frequency shiny app to explore how children's color input changes over development

³⁵⁷ Using childesr. We can analyze these learning trajectories using childesr by ³⁵⁸ breaking the process into five steps: (1) define our words of interest, (2) find the frequencies ³⁵⁹ with which children hear these words, (3) find the proportion of the *total words* children hear ³⁶⁰ that these frequencies account for, (4) aggregate across transcripts and children to determine ³⁶¹ the error in our estimates of these proportions, and (5) plot the results.

For this analysis, we will define our words of interest as the basic color words of English (except for gray, which children hear very rarely). We store these in the colors variable, and then use the get_types function from childesr to get the type frequency of each of these words in all of the corpora in CHILDES. For demonstration, we look only at the types produced by the speakers in each corpus tagged as Mother and Father. We also restrict ourselves to children from 1–5 years old (12–60 months), and look only at the North American English corpora.

```
colors <- c("black", "white", "red", "green", "yellow", "blue", "brown",
                      "orange", "pink", "purple")
color_counts <- get_types(collection = "Eng-NA",
                      role = c("Mother", "Father"),
                     age = c(12,60),
                      type = colors)
```

To normalize correctly (i.e., to ask what proportion of the input children hear consists of these color words), we need to know how many total words these children hear from their parents in these transcripts. To do this, we use the get_speaker_statistics function, which will return a total number of tokens (num_tokens) for each of these speakers.

Get the ids corresponding to all of the speakers we are interested in
parent_ids <- color_counts %>%
 distinct(collection_id, corpus_id, transcript_id, speaker_id)

```
# Find the total number of tokens produced by these speakers
parents <- parent_ids %>%
   left_join(get_speaker_statistics(collection = "Eng-NA")) %>%
   select(collection id, corpus id, transcript id, speaker id, num tokens)
```

We now join these two pieces of information together – how many times each speaker produced each color word, and how many total words they produced. We then group the data into 6-month age bins, and compute the proportion of tokens that comprise each color for each child in each 6-month bin. For comparability with the web app analysis, these proportions are converted to parts per million words.

Finally, we use non-parametric bootstrapping to estimate 95% confidence intervals for our estimates of the parts per million words of each color term with the tidyboot package.

```
count_estimates_with_error <- count_estimates %>%
  tidyboot::tidyboot_mean(parts) %>%
  left_join(graph_colors) %>%
  mutate(color = factor(color, levels = colors))
```

Figure 6 shows the results of these analyses: Input frequency varies substantially over the 1–5 year range for nearly every color word.



Figure 6. Color frequency as a function of age. Points represent means across transcripts, error bars represent 95% confidence intervals computed by nonparametric bootstrap

Gender has long been known to be an important factor for early vocabulary Gender. 382 growth, with girls learning more words earlier than boys (Huttenlocher, Haight, Bryk, 383 Seltzer, & Lyons, 1991). Parent-report data from ten languages suggest that female children 384 have larger vocabularies on average than male children in nearly every language (Eriksson et 385 al., 2012). Comparable cross-linguistic analysis of naturalistic production data has not been 386 conducted, however, and these differences are easy to explore using childesr. By pulling 387 data from the transcript by speaker table, a user has access to a set of derived linguistic 388 measures that are often used to evaluate a child's grammatical development. In this worked 389 example, we walk through a sample analysis that explores gender differences in early lexical 390 diversity. 391

First, we use the childesr function call get_speaker_statistics to pull data relating to the aforementioned derived measures for children and their transcripts. Note that ³⁹⁴ we exclusively select the children's production data, and exclude their caregivers' speech.

```
stats <- get_speaker_statistics(role = "Target_Child")</pre>
```

This childesr call retrieves data from all collections and corpora, including those languages for which there are very sparse data. In order to make any substantial inferences from our analysis, we begin by filtering the dataset to include only languages for which there are a large number of transcripts (> 500). We also restrict our analysis to children under the age of four years.

```
number_of_transcripts_threshold <- 500
max_age <- 4
included_languages <- stats %>%
filter(target_child_age < max_age * 365) %>%
count(language) %>%
filter(n > number_of_transcripts_threshold) %>%
pull(language)
```

Our transcript by speaker table contains multiple derived measures of lexical 400 diversity – here we use MTLD (McCarthy, 2005). MTLD is derived from the average length 401 of orthographic words that are above a pre-specified type-token ratio, making it more robust 402 to transcript length than simple TTR. We start by filtering to include only those children for 403 which a sex was defined in the transcript, who speak a language in our subset of languages 404 with a large number of transcripts, and who are in the appropriate age range. We then 405 compute an average MTLD score for each child at each age point by aggregating across 406 transcripts while keeping information about the child's sex and language. Note that one 407 child in particular, "Leo" in the eponymous German corpus, contained transcripts that were 408 a collection of his most complex utterances (as caregivers were instructed to record); this 409 child was excluded from the analysis. 410

The data contained in CHILDES is populated from a diverse array of studies reflecting 411 varying circumstances of data collection. This point is particularly salient in our gender 412 analysis due to potential non-independence issues that may emerge from the inclusion of 413 many transcripts from longitudinal studies. To account for non-independence, we fit a linear 414 mixed effects model with a *gender* * age (treated as a quadratic predictor) interaction as 415 fixed effects, child identity as a random intercept, and qender + aqe by language as a 416 random slope, the maximal converging random effects structure (Barr, Levy, Scheepers, & 417 Tily, 2013).² The plot below displays the average MTLD scores for various children at 418 different ages, split by gender, with a line corresponding to the prediction of our fit mixed 419 effects model. 420

This plot reveals a slight gender difference in linguistic productivity in young children, replicating the moderate female advantage found by Eriksson et al. (2012). The goal of this analysis was to showcase an example of using childesr to explore the CHILDES dataset. We also highlighted some of the potential pitfalls – sparsity and non-independence – that

² All code and analyses are available at https://github.com/langcog/childes-db-paper





emerge in working with a diverse set of corpora, many of which were collected in longitudinal
studies.

427 Teaching with childes-db

Teachers of courses on early language acquisition often In-class demonstrations. 428 want to illustrate the striking developmental changes in children's early language. One 429 method is to present static displays that show text from parent-child conversations extracted 430 from CHILDES or data visualizations of various metrics of production and input (e.g., MLU 431 or Frequency), but one challenge of such graphics is that they cannot be modified during a 432 lecture and thus rely on the instructor selecting examples that will be compelling to students. 433 In contrast, in-class demonstrations can be a powerful way to explain complex concepts 434 while increasing student engagement with the course materials. 435

436 Consider the following demonstration about children's first words. Diary studies and

large-scale studies using parent report show that children's first words tend to fall into a 437 fairly small number of categories: people, food, body parts, clothing, animals, vehicles, toys, 438 household objects, routines, and activities or states (Clark, 2009; Fenson et al., 1994; Tardif 439 et al., 2008). The key insight is that young children talk about what is going on around 440 them: people they see every day, e.g., toys and small household objects they can manipulate 441 or food they can control. To illustrate this point, an instructor could: 442 1. introduce the research question (e.g., What are the types of words that children first 443 produce?), 444 2. allow students to reflect or do a pair-and-share discussion with their neighbor, 445 3. show the trajectory of a single lexical item while explaining key parts of the 446 visualization (see Panel A of Figure 8), 447 4. elicit hypotheses from students about the kinds of words that children are likely to 448 produce, 449 5. make real-time queries to the web application to add students' suggestions and talk 450 through the updated plots (Panels B and C of Figure 8), and 451 6. finish by entering a pre-selected set of words that communicate the important 452 takeaway point (Panel D of Figure 8). 453 **Tutorials and programming assignments.** One goal for courses on applied 454 natural language processing (NLP) is for students to get hands-on experience using NLP 455 tools to analyze real-world language data. A primary challenge for the instructor is to decide 456 how much time should be spent teaching the requisite programming skills for accessing and 457

formatting language data, which are typically unstructured. One pedagogical strategy is to abstract away these details and avoid having students deal with obtaining data and formatting text. This approach shifts students' effort away from data cleaning and towards programming analyses that encourage the exploration and testing of interesting hypotheses. In particular, the childesr API provides instructors with an easy-to-learn method for giving students programmatic access to child language data.



Figure 8. Worked example of using the web applications for in-class teaching. Panels A-D show how an instructor could dynamically build a plot during a lecture to demonstrate a key concept in language acquisition.

For example, an instructor could create a programming assignment with the specific 464 goal of reproducing the key findings in the case studies presented above – color words or 465 gender. Depending on the students' knowledge of R, the instructor could decide how much of 466 the childesr starter code to provide before asking students to generate their own plots and 467 write-ups. The instructor could then easily compare students' code and plots to the expected 468 output to measure learning progress. In addition to specific programming assignments, the 469 instructor could use the childes-db and childesr workflow as a tool for facilitating 470 student research projects that are designed to address new research questions. 471

472

Conclusion

We have presented childes-db, a database formatted mirror of the CHILDES dataset. This database – together with the R API and web apps – facilitates the use of child language data. For teachers, students, and casual explorers, the web apps allow browsing and demonstration. For researchers interested in scripting more complex analyses, the API allows
them to abstract away from the details of the CHAT format and easily create reproducible
analyses of the data. We hope that these functionalities broaden the set of users who can
easily interact with CHILDES data, leading to future insights into the process of language
acquisition.

childes-db addresses a number of needs that have emerged in our own research and 481 teaching, but there are still a number of limitations that point the way to future 482 improvements. For example, childes-db currently operates only on transcript data, without 483 links to the underlying media files; in the future, adding such links may facilitate further 484 computational and manual analyses of phonology, prosody, social interaction, and other 485 phenomena by providing easy access to the video and audio data. Further, we have focused 486 on including the most common and widely-used tiers of CHAT annotation into the database 487 first, but our plan is eventually to include the full range of tiers. Finally, a wide range of 488 further interactive analyses could easily be added to the current suite of web apps. We invite 480 other researchers to join us in both suggesting and contributing new functionality as our 490 system grows and adapts to researchers' needs. 491

492

References

- Ambridge, B., Kidd, E., Rowland, C. F., & Theakston, A. L. (2015). The ubiquity of
 frequency effects in first language acquisition. *Journal of Child Language*, 42(2),
 239–273.
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
 68(3), 255–278.
- ⁴⁹⁹ Bååth, R. (2010). ChildFreq: An online tool to explore word frequencies in child language.
 ⁵⁰⁰ Lucs Minor, 16, 1–6.
- ⁵⁰¹ Bird, S., & Loper, E. (2004). NLTK: The natural language toolkit. In *Proceedings of the acl*

- 2004 on interactive poster and demonstration sessions (p. 31). Association for
 Computational Linguistics.
- ⁵⁰⁴ Brown, R. (1973). A first language: The early stages. Harvard U. Press.
- ⁵⁰⁵ Chang, F. (2017). The lucid language researcher's toolkit [computer software]. Retrieved ⁵⁰⁶ from http://www.lucid.ac.uk/resources/for-researchers/toolkit/
- ⁵⁰⁷ Clark, E. V. (2009). *First language acquisition*. Cambridge University Press.
- Demuth, K., Culbertson, J., & Alter, J. (2006). Word-minimality, epenthesis and coda
 licensing in the early acquisition of english. Language and Speech, 49(2), 137–173.
- ⁵¹⁰ Donoho, D. L. (2010). An invitation to reproducible computational research. *Biostatistics*, ⁵¹¹ 11(3), 385–388.
- Elman, J. L. (1993). Learning and development in neural networks: The importance of
 starting small. Cognition, 48(1), 71–99.
- 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(2), 326–343.
- Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing visual
 input in the first two years. *Cognition*, 152, 101–107.
- ⁵²⁰ Fenson, L., Dale, P. S., Reznick, J. S., Bates, E., Thal, D. J., Pethick, S. J., ... Stiles, J.
- (1994). Variability in early communicative development. Monographs of the Society
 for Research in Child Development, i–185.
- ⁵²³ Goldwater, S., Griffiths, T. L., & Johnson, M. (2009). A bayesian framework for word ⁵²⁴ segmentation: Exploring the effects of context. *Cognition*, 112(1), 21–54.
- Goodman, J. C., Dale, P. S., & Li, P. (2008). Does frequency count? Parental input and the acquisition of vocabulary. *Journal of Child Language*, 35(3), 515–531.
- ⁵²⁷ Huttenlocher, J., Haight, W., Bryk, A., Seltzer, M., & Lyons, T. (1991). Early vocabulary ⁵²⁸ growth: Relation to language input and gender. *Developmental Psychology*, 27(2),

236.529

541

- Kline, M. (2012). CLANtoR. http://github.com/mekline/CLANtoR/; GitHub. 530 doi:10.5281/zenodo.1196626 531
- MacWhinney, B. (2000). The childes project: The database (Vol. 2). Psychology Press. 532
- MacWhinney, B. (2014). The childes project: Tools for analyzing talk, volume ii: The 533 database. Psychology Press. 534
- MacWhinney, B., & Snow, C. (1985). The child language data exchange system. Journal of 535 Child Language, 12(2), 271–295. 536
- Malvern, D. D., & Richards, B. J. (1997). A new measure of lexical diversity. British Studies 537 in Applied Linguistics, 12, 58–71. 538
- Marcus, G. F., Pinker, S., Ullman, M., Hollander, M., Rosen, T. J., Xu, F., & Clahsen, H. 539 (1992). Overregularization in language acquisition. Monographs of the Society for 540 Research in Child Development, i–178.
- McCarthy, P. M. (2005). An assessment of the range and usefulness of lexical diversity 542 measures and the potential of the measure of textual, lexical diversity (mtld). 543
- Dissertation Abstracts International, 66, 12. 544
- McCarthy, P. M., & Jarvis, S. (2010). MTLD, vocd-d, and hd-d: A validation study of 545 sophisticated approaches to lexical diversity assessment. Behavior Research Methods, 546 42(2), 381-392.547
- Mevlan, S. C., Frank, M. C., Roy, B. C., & Levy, R. (2017). The emergence of an abstract 548 grammatical category in children's early speech. Psychological Science, 28(2), 549 181 - 192.550
- Miller, J. F., & Chapman, R. S. (1981). The relation between age and mean length of 551 utterance in morphemes. Journal of Speech, Language, and Hearing Research, 24(2), 552 154 - 161.553
- Montag, J. L., Jones, M. N., & Smith, L. B. (2015). The Words Children Hear: Picture 554
- Books and the Statistics for Language Learning. Psychological Science, 26(9), 555

- Norrman, G., & Bylund, E. (2015). The irreversibility of sensitive period effects in language
 development: Evidence from second language acquisition in international adoptees.
 Developmental Science.
- ⁵⁶⁰ R Core Team. (2017). R: A language and environment for statistical computing. Vienna,

561Austria: R Foundation for Statistical Computing. Retrieved from562https://www.R-project.org/

- Redington, M., Chater, N., & Finch, S. (1998). Distributional information: A powerful cue
 for acquiring syntactic categories. *Cognitive Science*, 22(4), 425–469.
- Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old
 infants. Science, 274 (5294), 1926–1928.
- ⁵⁶⁷ Snyder, W. (2007). Child language: The parametric approach. Oxford University Press.
- Song, J. Y., Shattuck-Hufnagel, S., & Demuth, K. (2015). Development of phonetic variants
 (allophones) in 2-year-olds learning american english: A study of alveolar stop/t,
 d/codas. Journal of Phonetics, 52, 152–169.
- Stodden, V., McNutt, M., Bailey, D. H., Deelman, E., Gil, Y., Hanson, B., ... Taufer, M.
 (2016). Enhancing reproducibility for computational methods. *Science*, 354 (6317),
 1240–1241.
- Tardif, T., Fletcher, P., Liang, W., Zhang, Z., Kaciroti, N., & Marchman, V. A. (2008).
 Baby's first 10 words. *Developmental Psychology*, 44 (4), 929.
- ⁵⁷⁶ Templin, M. (1957). Certain language skills in children: Their development and
- interrelationships (monograph series no. 26). Minneapolis: University of Minnesota,
 the Institute of Child Welfare.
- ⁵⁷⁹ Wagner, K., Dobkins, K., & Barner, D. (2013). Slow mapping: Color word learning as a ⁵⁸⁰ gradual inductive process. *Cognition*, 127(3), 307–317.
- ⁵⁸¹ Watkins, R. V., Kelly, D. J., Harbers, H. M., & Hollis, W. (1995). Measuring children's
- lexical diversity: Differentiating typical and impaired language learners. *Journal of*

⁵⁸³ Speech, Language, and Hearing Research, 38(6), 1349–1355.

- Wickham, H., & Grolemund, G. (2016). *R for data science: Import, tidy, transform, visualize, and model data.* "O'Reilly Media, Inc."
- Wickham, H., Francois, R., Henry, L., & Müller, K. (2017). Dplyr: A grammar of data
 manipulation. Retrieved from https://CRAN.R-project.org/package=dplyr
- Yang, C. (2013). Ontogeny and phylogeny of language. Proceedings of the National Academy
 of Sciences, 110(16), 6324–6327.
- Yurovsky, D., Wagner, K., Barner, D., & Frank, M. C. (2015). Signatures of domain-general
 categorization mechanisms in color word learning. In *CogSci*.