

## MIT Open Access Articles

*childes-db: A flexible and reproducible interface to the child language data exchange system*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**As Published:** <https://doi.org/10.3758/s13428-018-1176-7>

**Publisher:** Springer US

**Persistent URL:** <https://hdl.handle.net/1721.1/131922>

**Version:** Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

**Terms of Use:** Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.



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

3 Alessandro Sanchez\*<sup>1</sup>, Stephan C. Meylan\*<sup>2</sup>, Mika Braginsky<sup>3</sup>, Kyle E. MacDonald<sup>1</sup>,  
4 Daniel Yurovsky<sup>4</sup>, & Michael C. Frank<sup>1</sup>

5 <sup>1</sup> Stanford University

6 <sup>2</sup> University of California, Berkeley

7 <sup>3</sup> MIT

8 <sup>4</sup> University of Chicago

9 Author Note

10 Co-first authorship indicated with \*. Thanks to Brian MacWhinney for advice and  
11 guidance, and to Melissa Kline for her work on ClanToR, which formed a starting point for  
12 our work. This work is supported by a Jacobs Advanced Research Fellowship to MCF.

13 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](mailto:sanchez7@stanford.edu)

16

## Abstract

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

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

29 Word count: 2925

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

## 32 **Introduction**

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

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

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

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

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

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

80       Nonetheless, there are some outstanding challenges working with CHILDES, both for  
81 students and for advanced users. The CHILDES ecosystem uses a specialized file format  
82 (CHAT), which is stored as plain text but includes structured annotations grouped into  
83 parallel information "tiers" on separate lines. These tiers allow for a searchable plaintext

84 transcript of an utterance to be stored along with structured annotations of its  
85 phonological, morphological, or syntactic content. These files are usually analyzed using a  
86 command-line program (CLAN) that allows users to count word frequencies, compute  
87 statistics (e.g., mean length of utterance, or MLU), and execute complex searches against  
88 the data. While this system is flexible and powerful, mastering the CHAT codes and  
89 especially the CLAN tool with its many functions and flags can be daunting. These  
90 technical barriers decrease the ease of exploration by a novice researcher or in a classroom  
91 exercise.

92         On the opposite end of the spectrum, for data-oriented researchers who are interested  
93 in doing large-scale analyses of CHILDES, the current tools are also not ideal. CLAN  
94 software is an excellent tool for interactive exploration, but – as a free-standing application  
95 – it can be tricky to build into a processing pipeline written in Python or R. Thus,  
96 researchers who would like to ingest the entire corpus (or some large subset) into a  
97 computational analysis typically write their own parsers of the CHAT format to extract the  
98 subset of the data they would like to use (e.g., Kline, 2012; Meylan et al., 2017; Redington  
99 et al., 1998; Yang, 2013).

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

111 creating a full mirror of the data. These factors together lead to a lack of *computational*  
112 *reproducibility*, a major problem that keeps researchers from verifying or building on  
113 published research (Donoho, 2010; Stodden et al., 2016).

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

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

## 127 **Design and technical approach**

128 As described above, CHILDES is most often approached as a set of distinct CHAT  
129 files, which are then parsed by users, often using CLAN. In contrast to this parsing  
130 approach, which entails the sequential processing of strings, `childes-db` treats CHILDES  
131 as a set of linked tables, with records corresponding to intuitive abstractions such as words,  
132 utterances, and transcripts (see Kline, 2012 for an earlier example of deriving a singular  
133 tabular representation of a CHILDES transcript). Users of data analysis languages like R or  
134 Julia, libraries like Pandas, or those familiar with Structured Query Language (SQL) will  
135 be familiar with operations on tabular representations of data such as filtering (subsetting),  
136 sorting, aggregation (grouping), and joins (merges). These operations obviate the need for

137 users to consider the specifics of the CHAT representation – instead they simply request  
138 the entities they need for their research and allow the API to take care of the formatting  
139 details. We begin by orienting readers to the design of the system via a top-level  
140 description and motivation for the design of the database schema, then provide details on  
141 the database’s current technical implementation and the versioning scheme. Users  
142 primarily interested in accessing the database can skip these details and focus on access  
143 through the `childesr` API and the web apps.

#### 144 **Database format**

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

154 Many of the other tables in the database describe hierarchical collections built out of  
155 tokens – *utterance*, *transcript*, *corpus*, and *collection* – and store attributes appropriate for  
156 each level of description. Every entity includes attributes that link it to all higher-order  
157 collections, e.g., an utterance lists the transcript, corpus, and collection to which it belongs.  
158 An *utterance* contains one or more word tokens and includes fields such as the utterance  
159 type (e.g., *declarative*, *interrogative*, etc.), total number of tokens, and the total number of  
160 morphemes if the morphological structure is available in the original CHAT file. A  
161 *transcript* consists of one or more utterances and includes the date collected, the name of  
162 the target child, the age in days if defined, and the filename from CHILDES. A *corpus*



163 consists of one or more transcripts, corresponding to well-known collections like the Brown  
164 (Brown, 1973) or Providence (Demuth, Culbertson, & Alter, 2006) corpus. Finally, a  
165 *collection* is a superordinate collection of corpora generally corresponding to a geographic  
166 region, following the convention in CHILDES. Because every record can be linked to a  
167 top-level collection (generally corresponding to a language), each table includes data from  
168 all languages represented in CHILDES.

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

## 177 **Technical implementation**

178 `chilides-db` is stored as a MySQL database, an industry-standard, open-source  
179 relational database that can be accessed directly from a wide range of programming  
180 languages. The `chilides-db` project provides access to hosted, read-only databases on a  
181 publicly-accessible server for direct access and `chilidesr` (described below). The project  
182 also hosts compressed .sql exports for local installation. While the former is appropriate for  
183 most users, local installation can provide performance gains by allowing a user to access  
184 the database on their machine or on their local network, as well as allowing users to store  
185 derived information in the same database.

186 In order to import the CHILDES corpora into the MySQL schema described above, it  
187 must first be accurately parsed and subsequently vetted to ensure its integrity. We parse  
188 the XML (eXtensible Markup Language) release of CHILDES hosted by

189 [childes.talkbank.org](http://childes.talkbank.org) using the NLTK library in Python (Bird & Loper, 2004). Logic  
190 implemented in Python converts the linear, multi-tier parse into a tabular format  
191 appropriate for `childes-db`. This logic includes decisions that we review below regarding  
192 what information sources are captured in the current release of the database and which are  
193 left for future development.

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

203 **Versioning.** The content of CHILDES changes as additional corpora are added or  
204 transcriptions are updated; as of time of writing, these changes are not systematically  
205 tracked in a public repository.<sup>1</sup> To facilitate reproducibility of past analyses, we introduce  
206 a simple versioning system by adding a new complete parse of the current state of  
207 CHILDES every six months or as warranted by changes in CHILDES. By default, users  
208 interact with the most recent version of the database available. To support reproduction of  
209 results with previous versions of the database, we continue to host recent versions (up to  
210 the last three years / six versions) through our `childesr` API so that researchers can run  
211 analyses against specific historical versions of the database. For versions more than three  
212 years old, we host compressed `.sql` files that users may download and serve using a local

---

<sup>1</sup>Specific versions of the database, tracked using the version control system Git, can be obtained by emailing the maintainers of the CHILDES project. While tracking line-level changes with Git provides detailed information about what has changed, our method allows researchers to access the relevant version programmatically by simply adding an argument to a function call.

213 installation of MySQL server (for which we provide instructions).

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

- 218       • Sparsely annotated tiers, e.g. phonology (`%pho`) and situation (`%sit`)
- 219       • Media links
- 220       • Tone direction and stress
- 221       • Filled pauses
- 222       • Reformulations, word revision, and phrase revision, e.g. `<what did you>[//] how can`  
223       `you see it ?`
- 224       • paralinguistic material, e.g. `[=! cries]`

225 At present, `childes-db` focuses strictly on the contents of CHILDES, and does not include  
226 material in related TalkBank projects such as PhonBank, AphasiaBank, or DementiaBank.  
227 We will prioritize the addition of these information sources and others in response to  
228 community feedback.

## 229                               **Interfaces for Accessing `childes-db`**

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

### 232       **Interactive Web Apps**

233       The ability to easily browse and explore the CHILDES corpora is a cornerstone of the  
234 `childes-db` project. To this end we have created powerful yet easy-to-use interactive web  
235 applications that enable users to visualize various dimensions of the CHILDES corpus:  
236 frequency counts, mean lengths of utterance, type-token ratios, and more. All of this is

237 doable without the requirement of understanding command-line tools.<sup>2</sup>

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

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

259 **Derived Measures.** The syntactic complexity and lexical diversity of children’s  
260 speech are similarly critical metrics for acquisition researchers (Miller & Chapman, 1981;

---

<sup>2</sup>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, incremental sentence generation, and distributional word classification – our web apps focus on covering the same common tasks as CLAN, but yielding visualizations for the web browser.

261 Watkins, Kelly, Harbers, & Hollis, 1995). There are a number of well-established measures  
262 of children's speech that operationalize complexity and diversity, and have many  
263 applications in speech-language pathology (SLP), where measures outside of the normal  
264 range may be indicative of speech, language, or communication disorders.

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

- 270 • MLU-w (mean length of utterance in words),
- 271 • MLU-m (mean length of utterance in morphemes),
- 272 • TTR (type-token ratio, a measure of lexical diversity; Templin, 1957),
- 273 • MTLTD (measure of textual lexical diversity; Malvern & Richards, 1997),
- 274 • HD-D (lexical diversity via the hypergeometric distribution; McCarthy & Jarvis,  
275 2010)

276 As with the Frequency Counts app, a user may subset the data as they choose, compare  
277 measures between caregivers and children, and aggregate across children from different  
278 corpora.

279 **Population Viewer.** In many cases a researcher may want to view the statistics  
280 and properties of corpora (e.g., their size, number of utterances, number of tokens) before  
281 choosing a target corpus or set of corpora for an analysis. This web app is intended to  
282 provide a basic overview regarding the scale and temporal extent of various corpora in  
283 CHILDES, as well as give researchers insight into the aggregate characteristics of  
284 CHILDES. For example, examining the aggregate statistics reveals that coverage in  
285 CHILDES peaks at around 30 months.

## 286 **The `childesr` Package**

287       Although the interactive analysis tools described above cover some of the most  
288 common use cases of CHILDES data, researchers interested in more detailed and flexible  
289 analyses will want to interface directly with the data in `childes-db`. Making use of the R  
290 programming language (R Core Team, 2017), we provide the `childesr` package. R is an  
291 open-source, extensible statistical computing environment that is rapidly growing in  
292 popularity across fields and is increasing in use in child language research (Norrman &  
293 Bylund, 2015; e.g. Song, Shattuck-Hufnagel, & Demuth, 2015). The `childesr` package  
294 abstracts away the details of connecting to and querying the database. Users can take  
295 advantage of the tools developed in the popular `dplyr` package (Wickham, Francois, Henry,  
296 & Müller, 2017), which makes manipulating large datasets quick and easy. We describe the  
297 commands that the package provides and then give several worked examples of analyses  
298 using the package.

299       The `childesr` package is easily installed via CRAN, the comprehensive R archive  
300 network. To install, simply type: `install.packages("childesr")`. After installation,  
301 users have access to functions that can be used to retrieve tabular data from the database:

- 302       • `get_collections()` gives the names of available collections of corpora (“Eng-NA”,  
303       “Spanish”, etc.)
- 304       • `get_corpora()` gives the names of available corpora (“Brown”, “Clark”, etc.)
- 305       • `get_transcripts()` gives information on available transcripts (language, date, target  
306       child demographics)
- 307       • `get_participants()` gives information on transcript participants (name, role,  
308       demographics)
- 309       • `get_speaker_statistics()` gives summary statistics for each participant in each  
310       transcript (number of utterances, number of types, number of tokens, mean length of  
311       utterance)
- 312       • `get_utterances()` gives information on each utterance (glosses, stems, parts of

313 speech, utterance type, number of tokens, number of morphemes, speaker  
314 information, target child information)

- 315 • `get_types()` gives information on each type within each transcript (gloss, count,  
316 speaker information, target child information)
- 317 • `get_tokens()` gives information on each token (gloss, stem, part of speech, number  
318 of morphemes, speaker information, target child information)

319 Each of these functions take arguments that restrict the query to a particular subset of the  
320 data (e.g. by collection, by corpus, by speaker role, by target child age, etc.) and returns  
321 the output in the form of a table. All functions support the specification of the database  
322 version to use. For more detailed documentation, see the package repository  
323 (<http://github.com/langcog/childesr>).

## 324 Using `childes-db`: Worked Examples

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

### 329 Research applications

330 **Color frequency.** One common use of CHILDES is to estimate the frequency with  
331 which children hear different words. These frequency estimates are used both in the  
332 development of theory (e.g., frequent words are learned earlier; Goodman et al., 2008), and  
333 in the construction of age-appropriate experimental stimuli. One benefit of the `childes-db`  
334 interface is that it allows for easy analysis of how the frequencies of words change over  
335 development. Many of our theories in which children learn the structure of language from  
336 its statistical properties implicitly assume that these statistics are *stationary*,  
337 i.e. unchanging over development (e.g., Saffran, Aslin, & Newport, 1996). However a

338 number of recent analyses show that the frequencies with which infants encounter both  
339 linguistic and visual properties of their environment may change dramatically over  
340 development (Fausey, Jayaraman, & Smith, 2016), and these changing distributions may  
341 produce similarly dramatic changes in the ease or difficulty with which these regularities  
342 can be learned (Elman, 1993).

343 To demonstrate how one might discover such non-stationarity, we take as a case  
344 study the frequency with which children hear the color words of English (e.g. “blue”,  
345 “green”). Color words tend to be learned relatively late by children, potentially in part due  
346 to the abstractness of the meanings to which they refer (see Wagner, Dobkins, & Barner,  
347 2013). However, within the set of color words, the frequency with which these words are  
348 heard predicts a significant fraction of the variance in their order of acquisition (Yurovsky,  
349 Wagner, Barner, & Frank, 2015). But are these frequencies stationary – e.g. do children  
350 hear “blue” as often at 12 months as they do at 24 months? We answer this question in  
351 two ways – first using the web apps, and then using the `childesr` package.

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

362 From this figure, it seems likely that children hear “blue” more frequently early in  
363 development, but the trajectories of “red” and “green” are less clear. We also do not have a  
364 good sense of the errors of these measurements, are limited to just a few colors at a time



365 before the plot becomes too crowded, and cannot combine frequencies across speakers. To  
366 perform this analysis in a more compelling and complete way, a user can use the `childesr`  
367 interface.

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

373 For this analysis, we will define our words of interest as the basic color words of  
374 English (except for gray, which children hear very rarely). We store these in the `colors`  
375 variable, and then use the `get_types()` function from `childesr` to get the type frequency  
376 of each of these words in all of the corpora in CHILDES. All other functions are provided  
377 by base R or the `tidyverse` package. For demonstration, we look only at the types  
378 produced by the speakers in each corpus tagged as Mother and Father. We also restrict  
379 ourselves to children from 1–5 years old (12–60 months), and look only at the North  
380 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)
```

381 To normalize correctly (i.e., to ask what proportion of the input children hear consists  
382 of these color words), we need to know how many total words these children hear from their  
383 parents in these transcripts. To do this, we use the `get_speaker_statistics()` function,

384 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)
```

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

```
count_estimates <- color_counts %>%
  left_join(parents) %>%
  mutate(age_months = target_child_age,
         age_bin = as.integer(floor(age_months / 6) * 6),
         color = tolower(gloss)) %>%
  group_by(age_bin, color, target_child_id, transcript_id) %>%
  summarise(transcript_count = sum(count), transcript_num_tokens = sum(num_tokens)) %>%
  group_by(age_bin, color, target_child_id) %>%
  summarise(child_count = sum(transcript_count), child_num_tokens = sum(transcript_num_t
  mutate(parts = child_count / child_num_tokens * 1e6)
```

390 Finally, we use non-parametric bootstrapping to estimate 95% confidence intervals for  
 391 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))
```

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

394 **Gender.** Gender has long been known to be an important factor for early  
395 vocabulary growth, with girls learning more words earlier than boys (Huttenlocher, Haight,  
396 Bryk, Seltzer, & Lyons, 1991). Parent-report data from ten languages suggest that female  
397 children have larger vocabularies on average than male children in nearly every language  
398 (Eriksson et al., 2012). Comparable cross-linguistic analysis of naturalistic production data  
399 has not been conducted, however, and these differences are easy to explore using `childesr`.  
400 By pulling data from the `transcript_by_speaker` table, a user has access to a set of  
401 derived linguistic measures that are often used to evaluate a child’s grammatical  
402 development. In this worked example, we walk through a sample analysis that explores  
403 gender differences in early lexical diversity.

404 First, we use the `childesr` function call `get_speaker_statistics()` to pull data  
405 relating to the aforementioned derived measures for children and their transcripts. Note  
406 that we exclusively select the children’s production data, and exclude their caregivers’  
407 speech.

```
speaker_stats <- get_speaker_statistics(role = "Target_Child")
```

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

```

number_of_transcripts_threshold <- 500
max_age <- 4

included_languages <- speaker_stats %>%
  filter(target_child_age < max_age * 12) %>%
  count(language) %>%
  filter(n > number_of_transcripts_threshold) %>%
  pull(language)

```

413 Our `transcript_by_speaker` table contains multiple derived measures of lexical  
 414 diversity – here we use MTLT (McCarthy, 2005). MTLT is derived from the average  
 415 length of orthographic words that are above a pre-specified type-token ratio, making it  
 416 more robust to transcript length than simple TTR. We start by filtering to include only  
 417 those children for which a sex was defined in the transcript, who speak a language in our  
 418 subset of languages with a large number of transcripts, and who are in the appropriate age  
 419 range. We then compute an average MTLT score for each child at each age point by  
 420 aggregating across transcripts while keeping information about the child’s sex and  
 421 language. Note that one child in particular, “Leo” in the eponymous German corpus,  
 422 contained transcripts that were a collection of his most complex utterances (as caregivers  
 423 were instructed to record); this child was excluded from the analysis.

```

mtld_data <- speaker_stats %>%
  filter(!is.na(target_child_sex), target_child_name != "Leo",
         language %in% included_languages) %>%
  group_by(target_child_id, target_child_age, target_child_sex, language) %>%
  summarise(measure = mean(mtld)) %>%
  ungroup() %>%
  mutate(age_years = target_child_age / 12,

```

```

target_child_sex = factor(target_child_sex,
                           levels = c("male", "female")) %>%
filter(age_years < max_age)

```

424 The data contained in CHILDES is populated from a diverse array of studies  
 425 reflecting varying circumstances of data collection. This point is particularly salient in our  
 426 gender analysis due to potential non-independence issues that may emerge from the  
 427 inclusion of many transcripts from longitudinal studies. To account for non-independence,  
 428 we fit a linear mixed effects model with a *gender \* age* (treated as a quadratic predictor)  
 429 interaction as fixed effects, child identity as a random intercept, and *gender + age* by  
 430 language as a random slope, the maximal converging random effects structure (Barr, Levy,  
 431 Scheepers, & Tily, 2013).<sup>3</sup> The plot below displays the average MTLN scores for various  
 432 children at different ages, split by gender, with a line corresponding to the prediction of our  
 433 fit mixed effects model.

434 This plot reveals a slight gender difference in linguistic productivity in young children,  
 435 replicating the moderate female advantage found by Eriksson et al. (2012). The goal of  
 436 this analysis was to showcase an example of using `childesr` to explore the CHILDES  
 437 dataset. We also highlighted some of the potential pitfalls – sparsity and non-independence  
 438 – that emerge in working with a diverse set of corpora, many of which were collected in  
 439 longitudinal studies.

## 440 Teaching with childes-db

441 **In-class demonstrations.** Teachers of courses on early language acquisition often  
 442 want to illustrate the striking developmental changes in children’s early language. One  
 443 method is to present static displays that show text from parent-child conversations  
 444 extracted from CHILDES or data visualizations of various metrics of production and input  
 445 (e.g., MLU or Frequency), but one challenge of such graphics is that they cannot be

<sup>3</sup>All code and analyses are available at <https://github.com/langcog/childes-db-paper>

446 modified during a lecture and thus rely on the instructor selecting examples that will be  
447 compelling to students. In contrast, in-class demonstrations can be a powerful way to  
448 explain complex concepts while increasing student engagement with the course materials.

449 Consider the following demonstration about children’s first words. Diary studies and  
450 large-scale studies using parent report show that children’s first words tend to fall into a  
451 fairly small number of categories: people, food, body parts, clothing, animals, vehicles,  
452 toys, household objects, routines, and activities or states (Clark, 2009; Fenson et al., 1994;  
453 Tardif et al., 2008). The key insight is that young children talk about what is going on  
454 around them: people they see every day, e.g., toys and small household objects they can  
455 manipulate or food they can control. To illustrate this point, an instructor could:

- 456 1. introduce the research question (e.g., What are the types of words that children first  
457 produce?),
- 458 2. allow students to reflect or do a pair-and-share discussion with their neighbor,
- 459 3. show the trajectory of a single lexical item while explaining key parts of the  
460 visualization (see Panel A of Figure ??),
- 461 4. elicit hypotheses from students about the kinds of words that children are likely to  
462 produce,
- 463 5. make real-time queries to the web application to add students’ suggestions and talk  
464 through the updated plots (Panels B and C of Figure ??), and
- 465 6. finish by entering a pre-selected set of words that communicate the important  
466 takeaway point (Panel D of Figure ??) .

467 **Tutorials and programming assignments.** One goal for courses on applied  
468 natural language processing (NLP) is for students to get hands-on experience using NLP  
469 tools to analyze real-world language data. A primary challenge for the instructor is to  
470 decide how much time should be spent teaching the requisite programming skills for  
471 accessing and formatting language data, which are typically unstructured. One pedagogical  
472 strategy is to abstract away these details and avoid having students deal with obtaining

473 data and formatting text. This approach shifts students' effort away from data cleaning  
474 and towards programming analyses that encourage the exploration and testing of  
475 interesting hypotheses. In particular, the `childesr` API provides instructors with an  
476 easy-to-learn method for giving students programmatic access to child language data.

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

## 485 Conclusion

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

494 `childes-db` addresses a number of needs that have emerged in our own research and  
495 teaching, but there are still a number of limitations that point the way to future  
496 improvements. For example, `childes-db` currently operates only on transcript data,  
497 without links to the underlying media files; in the future, adding such links may facilitate  
498 further computational and manual analyses of phonology, prosody, social interaction, and

499 other phenomena by providing easy access to the video and audio data. Further, we have  
500 focused on including the most common and widely-used tiers of CHAT annotation into the  
501 database first, but our plan is eventually to include the full range of tiers. Finally, a wide  
502 range of further interactive analyses could easily be added to the current suite of web apps.  
503 We invite other researchers to join us in both suggesting and contributing new functionality  
504 as our system grows and adapts to researchers' needs.

## 505 References

- 506 Ambridge, B., Kidd, E., Rowland, C. F., & Theakston, A. L. (2015). The ubiquity of  
507 frequency effects in first language acquisition. *Journal of Child Language*, *42*(2),  
508 239–273.
- 509 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for  
510 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and*  
511 *Language*, *68*(3), 255–278.
- 512 Bååth, R. (2010). ChildFreq: An online tool to explore word frequencies in child language.  
513 *Lucs Minor*, *16*, 1–6.
- 514 Bird, S., & Loper, E. (2004). NLTK: The natural language toolkit. In *Proceedings of the*  
515 *acl 2004 on interactive poster and demonstration sessions* (p. 31). Association for  
516 Computational Linguistics.
- 517 Brown, R. (1973). *A first language: The early stages*. Harvard U. Press.
- 518 Chang, F. (2017). The lucid language researcher's toolkit [computer software]. Retrieved  
519 from <http://www.lucid.ac.uk/resources/for-researchers/toolkit/>
- 520 Clark, E. V. (2009). *First language acquisition*. Cambridge University Press.
- 521 Demuth, K., Culbertson, J., & Alter, J. (2006). Word-minimality, epenthesis and coda  
522 licensing in the early acquisition of english. *Language and Speech*, *49*(2), 137–173.
- 523 Donoho, D. L. (2010). An invitation to reproducible computational research. *Biostatistics*,



- 524           11(3), 385–388.
- 525 Elman, J. L. (1993). Learning and development in neural networks: The importance of  
526           starting small. *Cognition*, 48(1), 71–99.
- 527 Eriksson, M., Marschik, P. B., Tulviste, T., Almgren, M., Pérez Pereira, M., Wehberg, S.,  
528           ... Gallego, C. (2012). Differences between girls and boys in emerging language  
529           skills: Evidence from 10 language communities. *British Journal of Developmental*  
530           *Psychology*, 30(2), 326–343.
- 531 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing  
532           visual input in the first two years. *Cognition*, 152, 101–107.
- 533 Fenson, L., Dale, P. S., Reznick, J. S., Bates, E., Thal, D. J., Pethick, S. J., ... Stiles, J.  
534           (1994). Variability in early communicative development. *Monographs of the Society*  
535           *for Research in Child Development*, i–185.
- 536 Goldwater, S., Griffiths, T. L., & Johnson, M. (2009). A bayesian framework for word  
537           segmentation: Exploring the effects of context. *Cognition*, 112(1), 21–54.
- 538 Goodman, J. C., Dale, P. S., & Li, P. (2008). Does frequency count? Parental input and  
539           the acquisition of vocabulary. *Journal of Child Language*, 35(3), 515–531.
- 540 Huttenlocher, J., Haight, W., Bryk, A., Seltzer, M., & Lyons, T. (1991). Early vocabulary  
541           growth: Relation to language input and gender. *Developmental Psychology*, 27(2),  
542           236.
- 543 Kline, M. (2012). CLANtoR. <http://github.com/mekline/CLANtoR/>; GitHub.  
544           doi:10.5281/zenodo.1196626
- 545 MacWhinney, B. (2000). *The childes project: The database* (Vol. 2). Psychology Press.
- 546 MacWhinney, B. (2014). *The childes project: Tools for analyzing talk, volume ii: The*  
547           *database*. Psychology Press.
- 548 MacWhinney, B., & Snow, C. (1985). The child language data exchange system. *Journal*  
549           *of Child Language*, 12(2), 271–295.
- 550 Malvern, D. D., & Richards, B. J. (1997). A new measure of lexical diversity. *British*

- 551        *Studies in Applied Linguistics*, 12, 58–71.
- 552 Marcus, G. F., Pinker, S., Ullman, M., Hollander, M., Rosen, T. J., Xu, F., & Clahsen, H.  
553        (1992). Overregularization in language acquisition. *Monographs of the Society for*  
554        *Research in Child Development*, i–178.
- 555 McCarthy, P. M. (2005). An assessment of the range and usefulness of lexical diversity  
556        measures and the potential of the measure of textual, lexical diversity (mtld).  
557        *Dissertation Abstracts International*, 66, 12.
- 558 McCarthy, P. M., & Jarvis, S. (2010). MTLT, vocd-d, and hd-d: A validation study of  
559        sophisticated approaches to lexical diversity assessment. *Behavior Research*  
560        *Methods*, 42(2), 381–392.
- 561 Meylan, S. C., Frank, M. C., Roy, B. C., & Levy, R. (2017). The emergence of an abstract  
562        grammatical category in children’s early speech. *Psychological Science*, 28(2),  
563        181–192.
- 564 Miller, J. F., & Chapman, R. S. (1981). The relation between age and mean length of  
565        utterance in morphemes. *Journal of Speech, Language, and Hearing Research*, 24(2),  
566        154–161.
- 567 Montag, J. L., Jones, M. N., & Smith, L. B. (2015). The Words Children Hear: Picture  
568        Books and the Statistics for Language Learning. *Psychological Science*, 26(9),  
569        1489–1496.
- 570 Norrman, G., & Bylund, E. (2015). The irreversibility of sensitive period effects in  
571        language development: Evidence from second language acquisition in international  
572        adoptees. *Developmental Science*.
- 573 R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna,  
574        Austria: R Foundation for Statistical Computing. Retrieved from  
575        <https://www.R-project.org/>
- 576 Redington, M., Chater, N., & Finch, S. (1998). Distributional information: A powerful cue

- 577 for acquiring syntactic categories. *Cognitive Science*, 22(4), 425–469.
- 578 Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old  
579 infants. *Science*, 274(5294), 1926–1928.
- 580 Snyder, W. (2007). *Child language: The parametric approach*. Oxford University Press.
- 581 Song, J. Y., Shattuck-Hufnagel, S., & Demuth, K. (2015). Development of phonetic  
582 variants (allophones) in 2-year-olds learning american english: A study of alveolar  
583 stop/t, d/codas. *Journal of Phonetics*, 52, 152–169.
- 584 Stodden, V., McNutt, M., Bailey, D. H., Deelman, E., Gil, Y., Hanson, B., . . . Taufer, M.  
585 (2016). Enhancing reproducibility for computational methods. *Science*, 354(6317),  
586 1240–1241.
- 587 Tardif, T., Fletcher, P., Liang, W., Zhang, Z., Kaciroti, N., & Marchman, V. A. (2008).  
588 Baby’s first 10 words. *Developmental Psychology*, 44(4), 929.
- 589 Templin, M. (1957). Certain language skills in children: Their development and  
590 interrelationships (monograph series no. 26). *Minneapolis: University of Minnesota,*  
591 *the Institute of Child Welfare*.
- 592 Wagner, K., Dobkins, K., & Barner, D. (2013). Slow mapping: Color word learning as a  
593 gradual inductive process. *Cognition*, 127(3), 307–317.
- 594 Watkins, R. V., Kelly, D. J., Harbers, H. M., & Hollis, W. (1995). Measuring children’s  
595 lexical diversity: Differentiating typical and impaired language learners. *Journal of*  
596 *Speech, Language, and Hearing Research*, 38(6), 1349–1355.
- 597 Wickham, H., & Grolemund, G. (2016). *R for data science: Import, tidy, transform,*  
598 *visualize, and model data*. “O’Reilly Media, Inc.”
- 599 Wickham, H., Francois, R., Henry, L., & Müller, K. (2017). *Dplyr: A grammar of data*  
600 *manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- 601 Yang, C. (2013). Ontogeny and phylogeny of language. *Proceedings of the National*  
602 *Academy of Sciences*, 110(16), 6324–6327.
- 603 Yurovsky, D., Wagner, K., Barner, D., & Frank, M. C. (2015). Signatures of

604

domain-general categorization mechanisms in color word learning. In *CogSci*.