

1 Web-CDI: A system for online administration of the MacArthur-Bates Communicative
2 Development Inventories

³ Benjamin deMayo¹, Danielle Kellier², Mika Braginsky³, Christina Bergmann⁴, Cielke
⁴ Hendriks⁴, Caroline Rowland^{4,6}, Michael C. Frank⁵, & Virginia A. Marchman⁵

⁵ ¹ Princeton University

⁶ ² University of Pennsylvania

⁷ ³ Massachusetts Institute of Technology

⁴ Max Planck Institute for Psycholinguistics

⁵ Stanford University

¹⁰ ⁶ Radboud University

11

Abstract

12 Understanding the mechanisms that drive variation in children's language acquisition
13 requires large, population-representative datasets of children's word learning across
14 development. Parent report measures such as the MacArthur-Bates Communicative
15 Development Inventories (CDI) are commonly used to collect such data, but the traditional
16 paper-based forms make the curation of large datasets logistically challenging. Many CDI
17 datasets are thus gathered using convenience samples, often recruited from communities in
18 proximity to major research institutions. Here, we introduce Web-CDI, a web-based tool
19 which allows researchers to collect CDI data online. Web-CDI contains functionality to
20 collect and manage longitudinal data, share links to test administrations, and download
21 vocabulary scores. To date, over 3,500 valid Web-CDI administrations have been
22 completed. General trends found in past norming studies of the CDI (e.g., Feldman et al.,
23 2000) are present in data collected from Web-CDI: scores of children's productive
24 vocabulary grow with age, female children show a slightly faster rate of vocabulary growth,
25 and participants with higher levels of educational attainment report slightly higher
26 vocabulary production scores than those with lower levels of education attainment. We
27 also report results from an effort to oversample non-white, lower-education participants via
28 online recruitment ($N = 243$). These data showed similar age, gender and primary
29 caregiver education trends to the full Web-CDI sample, but this effort resulted in a high
30 exclusion rate. We conclude by discussing implications and challenges for the collection of
31 large, population-representative datasets.

32 *Keywords:* vocabulary development, parent report

33 Word count: X

34 Web-CDI: A system for online administration of the MacArthur-Bates Communicative
35 Development Inventories

36 Children vary tremendously in their vocabulary development (Fenson et al., 1994;
37 Frank, Braginsky, Yurovsky, & Marchman, 2021). Characterizing this variability is central
38 to understanding the mechanisms that drive early language acquisition, yet capturing this
39 variation in broad, diverse samples of children has been a significant challenge for cognitive
40 scientists for decades. The MacArthur-Bates Communicative Development Inventories
41 (MB-CDI, or CDI for short) are a set of commonly-used parent report instruments for
42 assessing vocabulary development in early childhood (Fenson et al., 2007) that were
43 introduced in part to create a cost-effective method for measuring variability across
44 individuals.

45 In this paper, we introduce a web-based tool, Web-CDI, which was developed to
46 address the need for collecting CDI data in an online format. Web-CDI allows researchers
47 to increase the convenience of CDI administration, further decrease costs associated with
48 data collection and entry (particularly with item-level data), and access participant
49 samples that have traditionally been difficult to reach in language development research.

50 Our purpose in this paper is twofold: first, we describe Web-CDI as a platform which
51 streamlines the process of collecting CDI data and collates the data in a way that
52 facilitates the creation of large-scale, multisite collaborative datasets. Second, we profile
53 usage of Web-CDI thus far, with a particular focus on broadening the reach of traditional
54 paper-based methods of collecting vocabulary development data.

55 **The Importance of Parent Report Data**

56 Gaining empirical traction on variation in children's early language requires reliable
57 and valid methods for measuring language abilities, especially in early childhood (8 to 30
58 months). Parent report is a mainstay in this domain. Parents' reports are based on their

59 daily experiences with the child, which are much more extensive than a researcher or
60 clinician can generally obtain. Moreover, they are less likely to be influenced by factors
61 that may mask a child's true ability in the laboratory or clinic (e.g., shyness). One widely
62 used set of parent-report instruments is the MacArthur-Bates Communicative Development
63 Inventories, originally designed for children learning American English (Fenson et al.,
64 2007). The American English CDIs come in several versions, two of which are Words &
65 Gestures (WG) for children 8 to 18 months, focusing on word comprehension and
66 production, as well as gesture use, and Words & Sentences (WS) for children 16 to 30
67 months, focusing on word production and sentence structure. Both the WG and WS
68 measures come in short forms with vocabulary checklists of approximately 90-100 words
69 (Fenson et al., 2000), and long forms, which contain vocabulary checklists of several
70 hundred items each. (An additional shorter form of the Web-CDI for children 30-37
71 months, CDI-III, also exists.) Together, the CDI instruments allow for a comprehensive
72 picture of milestones that characterize language development in early childhood. A
73 substantial body of evidence suggests that these instruments are both reliable and valid
74 (e.g., Fenson et al., 2007, 1994) leading to their widespread use in thousands of research
75 studies over the last few decades. Initial large-scale work to establish the normative
76 datasets for the American English CDI not only provided key benchmarks for determining
77 children's progress, but also documented the extensive individual differences that
78 characterize early language learning during this critical period of development (Bates et al.,
79 1994; Fenson et al., 1994). Understanding the origins and consequences of this variability
80 remains an important empirical and theoretical endeavor (e.g., Bates & Goodman, 2001;
81 Bornstein & Putnick, 2012; see also, Frank, Braginsky, Yurovsky, & Marchman, 2021).

82 The popularity of CDI instruments has remained strong over the years, leading to
83 extensions of the methodology to alternative formats and cross-language adaptations
84 (Fenson et al., 2000). Many teams around the world have adapted the CDI format to the
85 particular languages and communities (Dale, 2015). Importantly, these adaptations are not

simply translations of the original form but rather incorporate the specific features of different languages and cultures, since linguistic variability exists even among cultures that share a native language. As an example of this phenomenon, the word “Cheerios” is more common in the United States than it is in the United Kingdom; as a result, it might be expected that caregivers would report children’s knowledge of this word in the U.S. and not the U.K., even though English is the most common language in both countries. To date there are more than 100 adaptations for languages around the globe. Moreover, several research groups have developed shorter versions of the CDI forms by randomly sampling items from the full CDI and comparing participants’ responses to established norms (Mayor & Mani, 2019) or by developing computer adaptive tests (CATs) that use item response theory or Bayesian approaches to guide the selection of a smaller subset of items to which participants respond (Chai, Lo, & Mayor, 2020; Kachergis et al., 2021; Makransky, Dale, Havmose, & Bleses, 2016).

While the reliability and validity of the original CDI instruments is well-established for the American English versions of the forms and several others, most existing norming samples are skewed toward families with more years of formal education and away from non-white groups (Fenson et al., 2007). For example, representation in the American English norming samples is generally restricted to families living on the U.S. east and west coasts. Further, although paper survey administration is a time-tested method, increasingly researchers and participants would prefer to use an electronic method to administer and fill CDI forms, obviating the need to track (and sometimes mail) paper forms, and the need to key in hundreds of item-wise responses for each child.

Here, we report on our recent efforts to create and distribute a web-based version of the CDIs in order to address some of the limitations of the standard paper versions. Online administration of the CDI is not a novel innovation – a variety of research groups have created purpose-build platforms for administering the CDI in particular languages. For example, Kristoffersen et al. (2013) collected a large normative sample of Norwegian CDIs

113 using a custom online platform. Similarly, the Slovak adaptation of the CDI uses an online
114 administration format (Kapalková & Slanèová, 2007). And many groups have used general
115 purpose survey software such as Qualtrics and Survey Monkey to administer CDIs and
116 variants online (e.g., Caselli, Lieberman, & Pyers, 2020). The innovation of Web-CDI is to
117 provide a comprehensive researcher management interface for the administration of a wide
118 range of CDI forms, allowing researchers to manage longitudinal administrations, download
119 scores, and share links with parents easily, all while satisfying strong guarantees regarding
120 privacy and anonymity. Moreover, a key benefit of a unified data collection and storage
121 system such as Web-CDI is that data from disparate sources are combined into a single
122 repository. This substantially reduces the overhead efforts associated with bringing
123 together data collected by researchers across the world and allows for the analysis of large
124 comparative datasets with the power to detect general trends in vocabulary development
125 that may emerge across languages. Finally, due to an agreement between the CDI Advisory
126 Board and Brookes Publishing, the publisher of the print versions of the CDI suite,
127 Web-CDI is free of charge for those researchers who agree to contribute their data for the
128 renorming of the long form instruments.

129

Introducing Web-CDI

130 Web-CDI is a web-based platform for CDI administration and management.
131 Web-CDI allows researchers to communicate with families by sharing URLs (web links that
132 contain individual users' own administration of the Web-CDI) via email or social media,
133 facilitating access to families in areas distant from an academic institution and eliminating
134 costly mailings and laboratory visits. Web-CDI also standardizes electronic administration
135 and scoring of CDI forms across labs and institutions, making possible the aggregation of
136 CDI data for later reuse and comparison across administrations by different labs. Indeed,
137 researchers who use Web-CDI grant the CDI Advisory Board permission to access and
138 analyze the resulting data on an opt-out basis, providing a path towards continual

¹³⁹ improvement of CDI instruments. Since 2018, more than 3,500 CDIs have been collected
¹⁴⁰ by 15 research groups throughout the U.S. who are using Web-CDI, demonstrating the
¹⁴¹ potential for large-scale data collection and aggregation.

¹⁴² Below, we outline how Web-CDI is used. We begin by detailing the consent obtention
¹⁴³ process and participant experience. Second, we describe the interface that researchers use
¹⁴⁴ to collect data using Web-CDI, specifying a number of common use cases for the platform.

¹⁴⁵ **Participant interface**

¹⁴⁶ Participants can complete the Web-CDI on a variety of devices, including personal
¹⁴⁷ computers and tablets. Web-CDI can be also administered on a smartphone, although the
¹⁴⁸ experience is not as ideal for the user due to the length of the survey and the small screen.
¹⁴⁹ As Web-CDI moves in the future to incorporate more short forms and computer adaptive
¹⁵⁰ tests (CATs) formats (e.g., Chai, Lo, & Mayor, 2020; Makranksy, Dale, Havmose, & Bleses,
¹⁵¹ 2016; Mayor & Mani, 2019), smartphone-responsive design will become a priority.

¹⁵² When a participant clicks a URL shared by a researcher, they are directed to a
¹⁵³ website presenting their own personal administration of the Web-CDI. In some cases, they
¹⁵⁴ may be asked to read and accept a waiver of consent documentation, depending on
¹⁵⁵ whether the researcher has chosen to use that feature (see also Researcher Interface below).

¹⁵⁶ *Instructions.* After completing the first demographics page, participants are provided
¹⁵⁷ with detailed instructions that are appropriate for either the Words & Gestures or Words
¹⁵⁸ & Sentences version (see Figure 1 for an example of the instructions for how to determine
¹⁵⁹ whether the child “understands and says” a word, which is pertinent to both the Words &
¹⁶⁰ Gestures and Words & Sentences forms.). In addition, there are more detailed instructions
¹⁶¹ for completing the vocabulary checklist. Unlike the traditional paper versions, instructions
¹⁶² on how to properly choose responses are provided both in written and pictorial form. The
¹⁶³ pictorial instructions (Figure 1) aim to further increase caregivers’ understanding of how to

Instructions: v

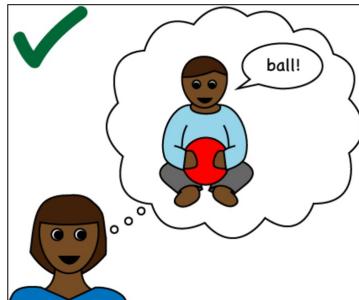
- This form can be filled anytime before the due date.
- It can also be saved at any time and resumed later by using the same link ([create bookmark](#)).
- After the form is submitted, it cannot be altered.
- The form also cannot be altered after the due date.
- Please use the navigation buttons below. Do not use the "back" and "forward" buttons on your browser.
- You can use the tab button and arrows keys to quickly navigate and answer questions.

Due date : Aug 8, 2017, 3:38 pm

Reach out to the Web-CDI Team!

Save

In this section, you will be asked about words that your child "understands and says." Your child "understands and says" a word on the list if they know what the word means AND they say it by themselves. Here are some examples. This assessment is for children of many ages. Your child may not be able to understand or say a lot of the words on the form. That is perfectly fine!



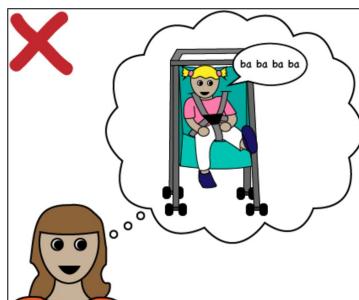
DO check the box if:

Your child says the word when trying to name an object or describe something that happened. You think s/he has a meaning for that word.

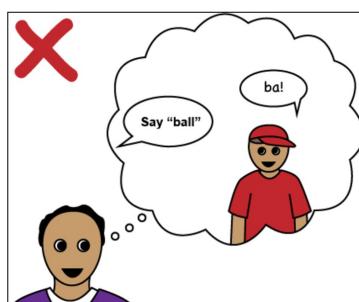


It's fine if your child can't say the whole word or says it his/her own "special" way. If you use another word in your family (e.g., Nana for Grandma), that is ok too!

DO check the box.



DON'T check the box if your child is just stringing sounds together. This is not a real word that means something.



DON'T ask your child to repeat the words on the list. This doesn't count!

Next Page >> 1/32**Save****Go back to Background Info**

Don't forget to save your progress whenever you can!

Figure 1. Pictorial instructions indicating how to mark whether a child *understands and says* a word, from the Web-CDI Words and Sentences instrument.

A**PART 1: Early Words****Vocabulary checklist**

The following is a list of typical words in young children's vocabularies. For words your child UNDERSTANDS but does not yet say, place a mark in the first column ("understands"). For words that your child both understands and also SAYS, place a mark in the second column ("understands and says"). You only need to mark one column. If your child uses a different pronunciation of a word (for example, "raffe" for "giraffe" or "sketti" for "spaghetti") or knows a different word that has a similar meaning as the word listed here (e.g., "nana" for "grandma"), go ahead and mark it. Remember, this is a "catalogue" of words that are used by many different children. Don't worry if your child knows only a few right now.

Hide/Show Instructions: ▾

1. Sound Effects And Animal Sounds

baa baa	<input type="checkbox"/> understands	<input type="checkbox"/> understands and says
choo choo	<input type="checkbox"/> understands	<input type="checkbox"/> understands and says
cockadoodledoo	<input type="checkbox"/> understands	<input type="checkbox"/> understands and says

B**PART 1: Words Children Use****A: Vocabulary Checklist**

Children understand many more words than they say. We are particularly interested in the words your child both understands and SAYS. Please go through the list and mark the words you have heard your child SAY on their own. If your child uses a different pronunciation of a word (for example, "raffe" instead of "giraffe" or "sketti" for "spaghetti") or says a different word that has a similar meaning as the word listed here (e.g., "nana" for "grandma"), go ahead and mark it. Remember that this is a "catalogue" of all the words that are used by many different children. Don't worry if your child only says a few of these right now.

Hide/Show Instructions: ▾

1. Sound Effects And Animal Sounds

<input type="checkbox"/> baa baa	<input type="checkbox"/> choo choo
<input type="checkbox"/> cockadoodledoo	<input type="checkbox"/> grr
<input type="checkbox"/> meow	<input type="checkbox"/> moo
<input type="checkbox"/> ouch	<input type="checkbox"/> quack quack
<input type="checkbox"/> uh oh	<input type="checkbox"/> vroom

Figure 2. (A) Sample items from the American English Words and Gestures form. (B) Sample items from the American English Words and Sentences form.

- 164 complete the checklist. For example, these instructions clarify that the child's
 165 understanding of a word requires them to have some understanding of the object that the
 166 word refers to or some aspect of the word's meaning. In addition, caregivers are reassured
 167 that "child-like" forms (e.g., "raff" for "giraffe") or family- or dialect-specific forms (e.g.,
 168 "nana" for "grandma") are acceptable evidence. Lastly, caregivers are reminded that the
 169 child should be able to produce the words "on their own" and that imitations are not
 170 acceptable. These general "rules of thumb" for completing the form should be familiar to
 171 researchers who are distributing the forms to caregivers so they can field any questions that
 172 may arise. While this is not possible for certain use-cases (e.g., social media recruitment),
 173 these instructions should ideally also be reviewed either in writing (e.g., via email) or

174 verbally (e.g., over the phone), so that these pictured instructions serve merely as a
175 reminder to caregivers when completing the form. Pictured instructions are available for
176 download on the MB-CDI website at <http://mb-cdi.stanford.edu/about.html>.

177 *Completing the instrument.* The majority of the participant's time is spent
178 completing the main sections of the instruments. As shown in Figure 2, on the American
179 English Words and Gestures form, the vocabulary checklist portion (396 items) asks
180 caregivers to indicate whether their child can "understand" or "understand and say" each
181 word; they can also indicate that their child neither understands nor says the word by
182 checking neither box. Additionally, gesture communication and other early milestones are
183 assessed. In the American English Words and Sentences form, the vocabulary checklist
184 (680 items) only asks caregivers to indicate which words their child "says." Additional
185 items assess children's production by requesting three of their longest sentences, as well as
186 morphological and syntactic development more broadly. All of these items are broken up
187 across multiple screens for easier navigation through the form.

188 At the completion of the form, a graph is displayed illustrating the proportion of
189 words from each semantic category that the child currently produces or understands.
190 Participants can select to download their own responses. In addition, data from the
191 norming studies are used to estimate the 'hardest' (i.e., most advanced based on previous
192 work on age of acquisition of individual words, Frank, Braginsky, Yurovsky, and Marchman
193 (2021)) word that the child currently understands or produces. This feedback to caregivers
194 is intended to provide caregivers with a fun "thank you" and intentionally avoids any
195 information which frames their child's progress relative to other children or any normative
196 standard, so as to not give the impression that the Web-CDI is a clinical assessment of the
197 child's development. To further underscore this point, the closing page reminds caregivers
198 that their participation does not constitute a clinical evaluation and that they should
199 contact their pediatrician or primary care physician if they have any concerns about their
200 child's development.

201 **Researcher interface**

202 One of the main goals of Web-CDI is to provide a unified CDI platform to the child
203 language research community. To that end, researchers request an account by contacting a
204 member of the CDI Advisory Board. Once the request is granted, they can design and
205 distribute studies. One rationale for this personalized registration process is that we ask
206 that researchers allow fully anonymized data from their participants to be shared with the
207 CDI Advisory Board, so that it can be added to Wordbank
208 [<http://wordbank.stanford.edu/>; Frank et al. (2017)] and shared with the broader research
209 community. However, if particular participants indicate in the consent process that they do
210 not want their data to be shared more broadly, then researchers can indicate this in the
211 Web-CDI dashboard to prevent data from specific administrations being contributed to any
212 analyses conducted by the CDI Advisory Board and/or Wordbank. Data currently in
213 Web-CDI, which have not yet been added to the Wordbank repository, will be vetted before
214 being added to ensure that all data being added to Wordbank from Web-CDI are drawn
215 from families with typically-developing children who meet similar inclusion criteria to the
216 ones we describe below in the *Dataset 1* section. Additionally, date of form completion will
217 be preserved when adding Web-CDI data into Wordbank, so that researchers can choose to
218 filter out data that may be affected by the particular point in time at which they were
219 collected (for example, the COVID-19 pandemic, Kartushina et al., 2021).

220 A study in the context of the Web-CDI system is a set of individual administrations
221 created by a researcher that share certain specifications. Table A1 in the Appendix gives
222 an overview of the customizable features that are available at the study level in Web-CDI.
223 These features are set when creating a study using the “Create Study” tool, and most of
224 the features can be updated continuously during data collection using the “Update Study”
225 tool. While some of these features are only relevant to specific use cases (e.g., longitudinal
226 research and social media data collection, described below), others are relevant to all

227 researchers using Web-CDI.

228 There are currently several CDI forms available for distribution on Web-CDI,
229 including the English WG and WS forms and forms in other languages (see Cross-linguistic
230 research, below). When creating a study, researchers choose one of the forms that they
231 would like to distribute to participants; only one can be used in a given study. Researchers
232 who wish to send multiple forms to participants simultaneously (e.g., those conducting
233 multilingual research) should create multiple studies, each with a single instrument
234 associated with it.

235 Researchers can download participant data in two formats. Both formatting options
236 output a comma-separated values file with one row per participant; the full data option
237 includes participant-by-item responses, and allows researchers to explore item-level trends,
238 while the summary data option omits item-level data and only provides summary scores
239 and normative information, including total number of words understood/produced and
240 percentile scores by age in months and gender. Percentile scores based are calculated to a
241 single percentile resolution using norms from Fenson et al. (2007).

242 Below, we outline several possible use cases of Web-CDI, as well the features which
243 may facilitate them from a researcher's perspective.

244 *Individual recruitment.* A first possible workflow using Web-CDI is to send unique
245 study URLs to individual participants. Researchers do so by entering numerical participant
246 IDs or by auto-generating a specified quantity of participant IDs, each with its own unique
247 study URL, using the “Add Participants” tool in the researcher dashboard. New
248 participants can be added on a continual basis so that researchers can adjust the sample
249 size of their study during data collection. Unique links generated for individual participants
250 expire, by default, 14 days after creation, though the number of days before link expiration
251 is adjustable, which may be an important consideration for some researchers depending on
252 their participant populations and specific project timelines. Workflows that involve

253 generating unique links are most suitable for studies which pair the CDI with other
254 measures, or when researchers contact specific participants from an existing database.

255 *Longitudinal studies.* Web-CDI also facilitates longitudinal study designs in which
256 each participant completes multiple administrations. Researchers wishing to design
257 longitudinal studies can do so by entering a list of meaningful participant IDs using the
258 “Add Participants” tool in the researcher dashboard. If a specific participant ID is added
259 multiple times, Web-CDI will automatically create multiple unique study URLs in the
260 study dashboard that have that ID. In addition, when creating studies, researchers can
261 select whether they would like the demographics information, vocabulary checklist, or no
262 sections at all to be pre-filled when a participant fills out a repeat administration of the
263 instrument. Unless researchers are interested in cumulative vocabulary counts, it is
264 strongly recommended that they do not use the option to pre-fill the vocabulary checklist
265 portion of the instrument in longitudinal administrations as caregivers should complete the
266 instrument at each time point independently. In the case that researchers do choose this
267 option, this is recorded in the Web-CDI database so that, when the data are added to
268 WordBank, researchers can choose to filter out any pre-filled questionnaires.

269 *Social media and survey vendors.* Web-CDI contains several features designed to
270 facilitate data collection from social media recruitment or through third-party
271 crowd-sourcing applications and vendors (e.g., Amazon Mechanical Turk, Prolific). First,
272 rather than creating unique survey links for each participant, researchers can also use a
273 single, anonymous link. When a participant clicks the anonymous link, a new
274 administration with a unique subject ID is created in the study dashboard. Additionally,
275 Web-CDI studies have several customizable features that are geared towards anonymous
276 online data collection. For example, researchers can adjust the minimum amount of time a
277 participant must take to fill out the survey before they are able to submit; with a longer
278 minimum time to completion, researchers can encourage a more thorough completion of the
279 survey. This feature is typically most relevant in research designs in which participants are

280 not vetted by the researcher or those in which there is no direct communication between
281 participants and researchers, as might be the case when recruiting respondents on social
282 media. Responses collected via personal communication with participants show low rates of
283 too-fast responding, mostly removing the need for the minimum time feature. Even in the
284 case of anonymous data collection, however, it is recommended that researchers not raise
285 the minimum completion time higher than 6 minutes, since some caregivers of very young
286 children may theoretically be able to proceed through the measure quickly if their child is
287 not yet verbal. Aside from the minimum time feature, researchers can ask participants to
288 verify that their information is accurate by checking a box at the end of the survey, and
289 can opt to include certain demographic questions at both the beginning and end of the
290 survey, using response consistency on these redundant items as a check of data quality.

291 *Paid participation.* If researchers choose to compensate participants directly through
292 the Web-CDI interface, Web-CDI has built-in functionality to distribute redeemable gift
293 codes when a participant reaches the end of the survey. Web-CDI contains several features
294 to facilitate integration with third-party crowdsourcing applications and survey vendors
295 should they choose to handle participant compensation through another platform. For
296 example, when creating studies, researchers can enter a URL to which participants are
297 redirected when they reach the end of the survey. Researchers using the behavioral
298 research platform Prolific can configure their study to collect participants' unique Prolific
299 IDs and pre-fill them in the survey.

300 *Cross-linguistic research.* Web-CDI forms are currently available in English (U.S.
301 American and Canadian), Spanish, French (Quebecois), Hebrew, Dutch and Korean. We
302 are looking to add more language forms to the tool, as the paper version of the forms has
303 been adapted into more than 100 different languages and dialects, and further ongoing
304 adaptations have been approved by the MB-CDI board
305 (<http://mb-cdi.stanford.edu/adaptations>).

306 System Design

307 Web-CDI is constructed using open-source software. All of the vocabulary data
308 collected in Web-CDI are stored in a standard MySQL relational database, managed using
309 Django and Python and hosted either by Amazon Web Services or by a European Union
310 (GDPR) compliant server (see below). Individual researchers can download data from their
311 studies through the researcher interface, and Web-CDI administrators have access to the
312 entire aggregate set of data from all studies run with Web-CDI. Website code is available in
313 a GitHub repository at <https://github.com/langcog/web-cdi>, where interested users can
314 browse, make contributions, and request technical fixes.

315 Data Privacy and GDPR Compliance

316 Web-CDI is designed to be compliant with stringent human subjects privacy
317 protections across the world. First, for U.S. users, we have designed Web-CDI based on the
318 United States Department of Health and Human Services “Safe Harbor” Standard for
319 collecting protected health information as defined by the Health Insurance Portability and
320 Accountability Act (HIPAA). In particular, participant names are never collected, birth
321 dates are used to calculate age in months (with no decimal information) but never stored,
322 and geographic zip codes are trimmed to the first 3 digits. Because of the architecture of
323 the site, even though participants enter zip codes and dates of birth, these are never
324 transmitted in full to the Web-CDI server. Since no identifying information is being
325 collected by the Web-CDI system, this feature ensures that Web-CDI can be used by
326 United States labs without a separate Institutional Review Board agreement between
327 users’ labs and Web-CDI (though of course researchers using the site will need Institutional
328 Review Board approval of their own research projects).¹

¹ Issues of de-identification and re-identifiability are complex and ever changing. In particular, compliance with DHHS “safe harbor” standards does not in fact fully guarantee the impossibility of statistical

329 In the European Union (EU), research data collection and storage is governed by the

330 Generalized Data Protection Regulation (GDPR) and its local instantiation in the legal

331 system of the member states. Some of the questions on the demographic form contain

332 information that may be considered sensitive (e.g., information about children's

333 developmental disorders), and in some cases, the possibility of linking this sensitive

334 information to participant IDs exists, particularly when researchers draw on local databases

335 that contain full names and addresses for recruitment and contacting. As a result, issues

336 regarding GDPR compliance arise when transferring data outside the EU, namely to

337 Amazon Web Services servers housed in the United States. Following GDPR regulations,

338 these issues would make a data sharing agreement between data collectors and Amazon

339 Web Services necessary. In addition, all administrators who can access the collected data

340 would have to enter such an agreement, which needs updating whenever personnel changes

341 occur. To overcome these hurdles, and in consultation with data protection officers, we

342 opted to leverage the local technical expertise and infrastructure to set up a sister site

343 housed on GDPR-compliant servers, currently available at <http://webcdi.mpi.nl>. This site

344 is updated synchronously with the main Web-CDI website to ensure a consistent user

345 experience and access to the latest features and improvements. This site has been used in

346 135 successful administrations so far and is the main data collection tool for an ongoing

347 norming study in the Netherlands. We are further actively advertising the option to use

348 the European site to other labs who are following GDPR guidelines and are planning

349 adaptations to multiple European languages, where copyright allows.

350 Current data collection

351 We now turn to an overview of the data collected thus far using Web-CDI. First, we

352 examine the full sample of all of the Web-CDI administrations collected as of autumn 2020

re-identification in some cases and if potential users have questions, we encourage them to consult with an Institutional Review Board.

353 (Dataset 1); we then focus in on a specific subset of Dataset 1 which is comprised of data
 354 from recent efforts to oversample non-white, less highly-educated U.S. participants
 355 (Dataset 2). Across both datasets, we show that general trends from prior research on
 356 vocabulary development are replicated using Web-CDI. Based on this work to date, we
 357 then discuss the potential for using Web-CDI to collect vocabulary development data from
 358 diverse communities online.

359 **Dataset 1: Full Current Web-CDI Usage**

Table 1

Exclusions from Dataset 1: full Web-CDI sample

Exclusion	WG	% of full	WS	% of full
	exclusions	WG sample	exclusions	WS sample
		excluded		excluded
Not first administration	163	5.68%	444	12.35%
Premature or low birthweight	37	1.29%	67	1.86%
Multilingual exposure	449	15.66%	492	13.69%
Illnesses/Vision/Hearing	191	6.66%	203	5.65%
Out of age range	88	3.07%	200	5.56%
Completed survey too quickly	319	11.12%	274	7.62%
System error in word tabulation	1	0.03%	4	0.11%
Total exclusions	1248	44%	1684	47%

360 In this section, we provide some preliminary analyses of Dataset 1, which consists of
 361 the full sample of American English Web-CDI administrations collected before autumn
 362 2020. At time of writing, researchers from 15 universities in the United States have
 363 collected over 5,000 administrations of the American English CDI using Web-CDI since it
 364 was launched in late 2017, with 2,868 administrations of the WG form before exclusions
 365 and 3,594 administrations of the WS form before exclusions. We excluded participants

366 from the subsequent analyses based on the following set of stringent criteria designed for
367 the creation of future normative datasets. We excluded participants if it was not their first
368 administration of the survey; if they were born prematurely or had a birthweight under 5.5
369 lbs (< 2.5 kg); reported more than 16 hours of exposure to a language other than English
370 per week on average (amounting to approximately > 10% of time during a week that a
371 child hears another language than English); had serious vision impairments, hearing
372 deficits or other developmental disorders or medical issues²; were outside of the correct age
373 range for the survey; or spent less time on the survey than a pre-specified timing cutoff.
374 Timing cutoffs were determined by selecting two studies within Dataset 1 that, upon a
375 visual inspection, appeared to contain high-quality responses (i.e., did not contain a
376 disproportionate number of extremely quick responders), and using these to estimate the
377 5th percentile of completion time by the child's age in months with a quantile regression.
378 Thus, for each age on the WG and WS measures, we obtained an estimate of the 5th
379 percentile of completion time and used this estimate as the shortest amount of time
380 participants could spend on the Web-CDI without being excluded from our analyses here.

381 The exclusion criteria we used were designed to be generally comparable with those
382 used in Fenson et al. (2007), who adopted stringent criteria to establish vocabulary norms
383 that reflect typically developing children's vocabulary trajectories. A complete breakdown
384 of the number of participants excluded on each criterion is in Table 1. Of the completed
385 WG forms, 1,248 were excluded, leading to a final WG sample size of 1,620 administrations,
386 and 1,694 WS administrations were excluded, leading to a final WS sample size of 1,900.

387 **Demographic distribution and exclusions.** Figure 3 shows the distribution of
388 participant ethnicities in Dataset 1 as compared with previously reported numbers in the
389 published norming study of the paper-based CDI form by Fenson et al. (2007). Several
390 issues pertaining to sample representativeness are appreciable. First, as shown in Figure

² Exclusions on the basis of child health were decided on a case-by-case basis by author V.M. in consultation with Philip Dale, Donna Thal, and Larry Fenson.

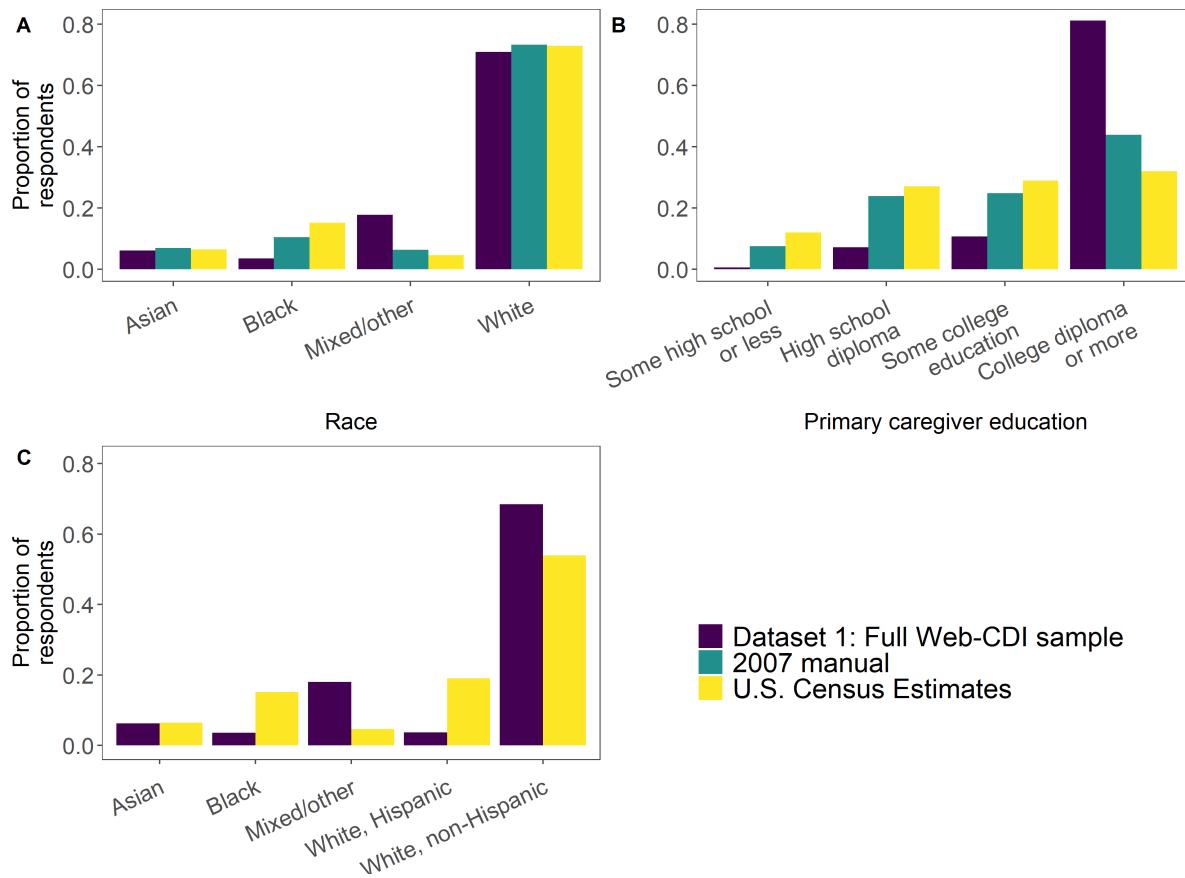


Figure 3. Top row: Proportion of respondents plotted by child race (A) and educational level of primary caregiver (B) from full Web-CDI sample (Dataset 1) to date ($N = 3,520$), compared with norming sample demographics from Fenson (2007) and U.S. Census data (American Community Survey, 2019; National Center for Education Statistics, 2019). Bottom row (C): Participant breakdown by race in Dataset 1 as compared with U.S. Census data, splitting white participants into those who are Hispanic and those are not.

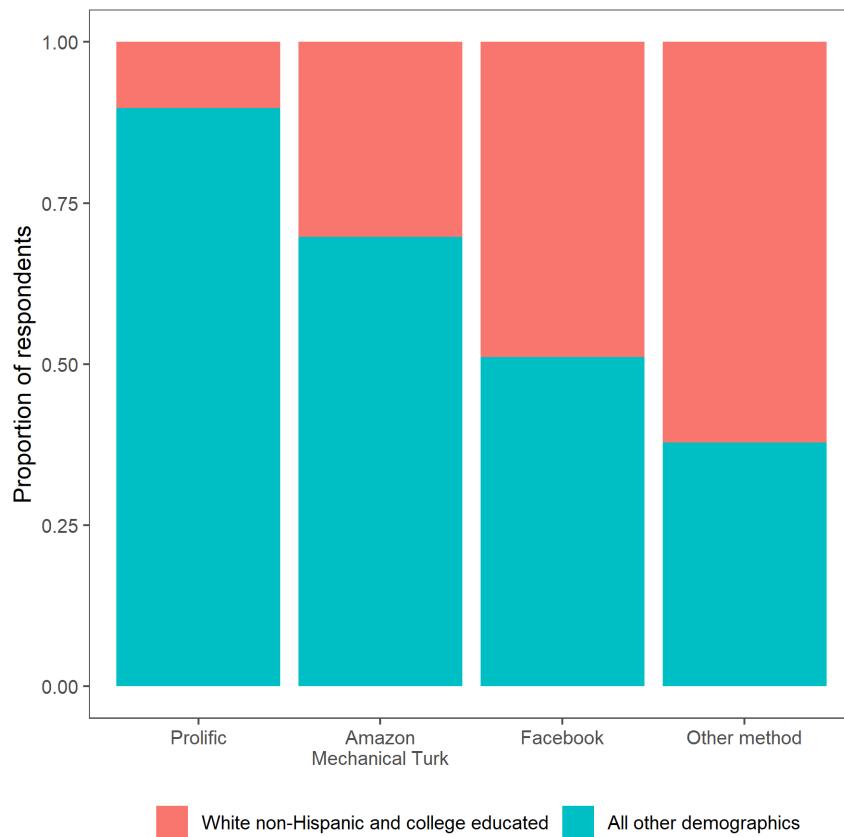


Figure 4. Proportion of participants from Dataset 1 who were white, college educated and not Hispanic, plotted by recruitment method.

391 3A, white participants comprised nearly three quarters of Dataset 1, which is comparable
 392 to U.S Census estimates in 2019 of U.S. residents between the ages of 15 and 34 in 2019;
 393 however, Figure 3C shows that, compared with U.S. Census estimates, many more white
 394 participants in Dataset 1 were non-Hispanic than is true of the U.S. population in general,
 395 indicating that Web-CDI is significantly oversampling white, non-Hispanic individuals (the
 396 breakdown of white participants into Hispanic and non-Hispanic is not reported in the
 397 2007 norms). Moreover, few participants identified as Hispanic/Latinx: 6.4% of WG
 398 participants and 5.2% of WS participants reported Hispanic or Latinx heritage. The low
 399 percentage of Hispanic/Latinx participants was due in part to our exclusion of children
 400 with substantial exposure to languages other than English: before exclusions, 8.4% of WG

401 participants were Hispanic/Latinx, and 8.1% of WS participants were Hispanic/Latinx.
402 Finally, representation of Black participants is generally lower in Dataset 1 (3.5%) than in
403 the 2007 norms (10.5%), which is in turn lower than U.S. Census estimates (15.2%). This
404 indicates that both Web-CDI data and existing norming samples tend to substantially
405 underrepresent Black participants.

406 Participants' educational attainment level, as measured by the primary caregiver's
407 highest educational level reached³, was similarly skewed. In Dataset 1, 81.2% of responses
408 came from families with college-educated primary caregivers compared to 43.8% from the
409 same group in the 2007 norms and 32.0% (Figure 3). Furthermore, less than 1% of
410 participants report a primary caregiver education level less than a high school degree,
411 compared to 7% from the same group in the 2007 norms.

412 The overrepresentation of white, non-Hispanic Americans and those with high levels
413 of education attainment points to a general challenge encountered in vocabulary
414 development research, which we return to when we detail our efforts to recruit more diverse
415 participants. Figure 4 shows that, of the recruitment methods used in Dataset 1, the
416 studies conducted using the platform Prolific (which we detail in the *Dataset 2* section)
417 contributed the least to the high proportion of white, non-Hispanic, college educated
418 participants. Respondents not known to be recruited through an online channel or
419 crowdsourcing platform (labeled "Other method" in Figure 4) showed the most
420 overrepresentation of white, college educated participants, suggesting that reliance on
421 university convenience samples may be driving the demographic skewness of Dataset 1
422 most acutely.

³ Maternal education level is a common measure of family socioeconomic status; we probe *primary caregiver* education level here to accommodate family structures in which child-rearing may not primarily be the responsibility of the child's mother, but we expect that in the vast majority of cases this corresponds to the child's mother.

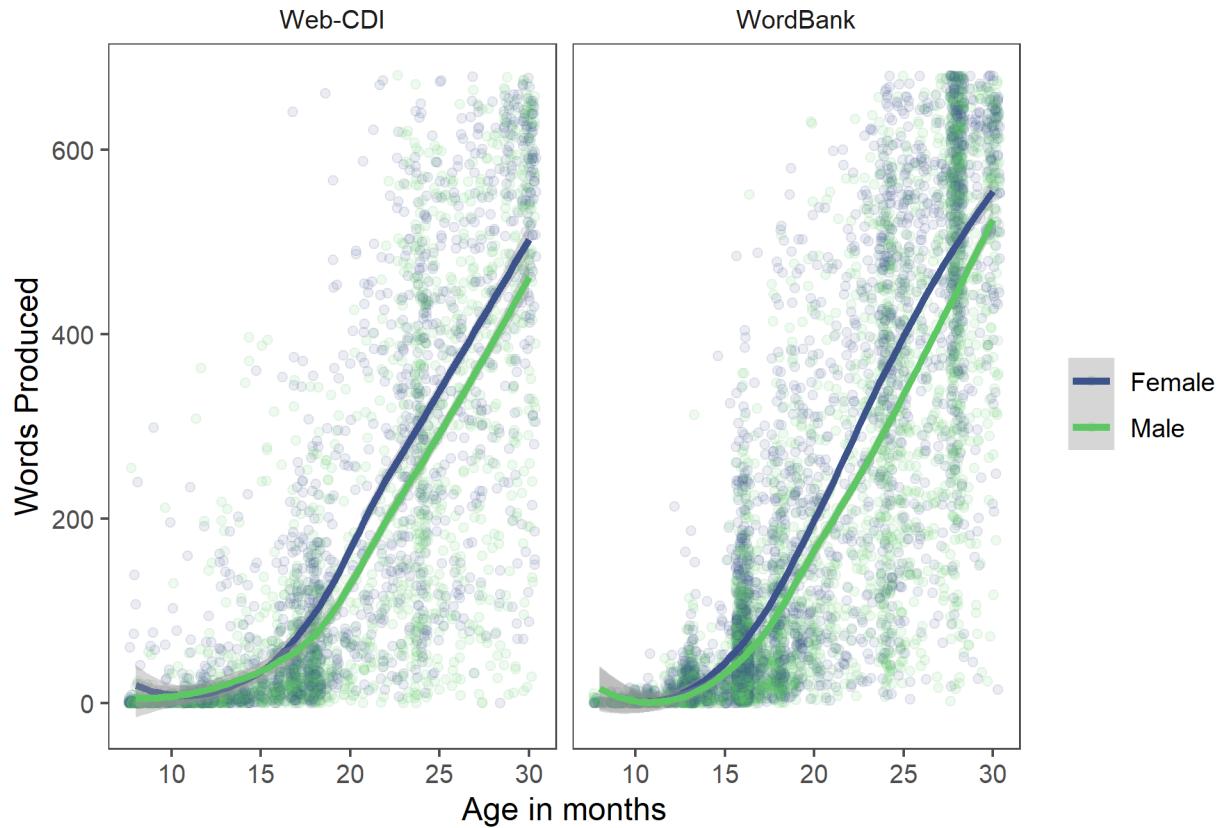


Figure 5. Individual children's vocabulary production scores plotted by children's age and gender (both WG and WS). Left panel: Dataset 1 (full sample of Web-CDI administrations, N = 3,510, with 1,673 girls). Right panel: American English CDI administrations in the WordBank repository (Frank et al., 2021), including only those administrations for which the child's gender was available (N = 6,486, with 3,146 girls). Lines are locally weighted regressions (LOESS) with associated 95% confidence intervals. Children with a different or no reported gender (N = 10) are omitted here.

423 **Results: Dataset 1.** Although the CDI instruments include survey items intended
424 to measure constructs other than vocabulary size, such as gesture, sentence production and
425 grammar, we focus exclusively on the vocabulary measures here. We also visualize key
426 analyses from Dataset 1 alongside the analogous analyses on the American English CDI
427 administrations from the WordBank repository (Frank, Braginsky, Yurovsky, & Marchman,
428 2021) that include the relevant demographic information needed to provide a comparison
429 dataset of traditional paper-and-pencil forms. Across both the WG and WS measures,
430 Dataset 1 shows greater reported vocabulary comprehension and production for older
431 children. Moreover, data from both the WG and WS measures in Dataset 1 replicate a
432 subtle but reliable pattern such that female children tend to have slightly larger vocabulary
433 scores than male children across the period of childhood assessed in the CDI forms (Frank,
434 Braginsky, Yurovsky, & Marchman, 2021), though in these data this difference does not
435 appear until around 18 months (Figure 5).

436 On the WG form, respondents' reports of children's vocabulary comprehension and
437 production both increased with children's age (Figure 6). We replicate overall patterns
438 found by Feldman et al. (2000) in that, on both the "Words Understood" and "Words
439 Produced" measures, vocabulary scores were slightly negatively correlated with primary
440 caregivers' education level, such that those caregivers without any college education
441 reported higher vocabulary scores on both scales; on the word comprehension scale, this
442 was particularly the case for the youngest infants in the sample. A linear regression model
443 with robust standard errors predicting comprehension scores with children's age and
444 primary caregivers' education level (binned into categories of "High school diploma or less,"
445 "Some college education" and "College diploma or more"⁴) as predictors shows main effects
446 of both age ($\beta = 20.05, p < 0.001$) and caregiver primary education ($\beta_{highschool} = 21.86, p$
447 = 0.05). Similarly, a linear regression model with robust standard errors predicting

⁴ "High school diploma or less" corresponds to 12 or fewer years of education; "Some college" corresponds to 13 - 15 years of education; "College diploma or more" refers to 16 or more years of education.

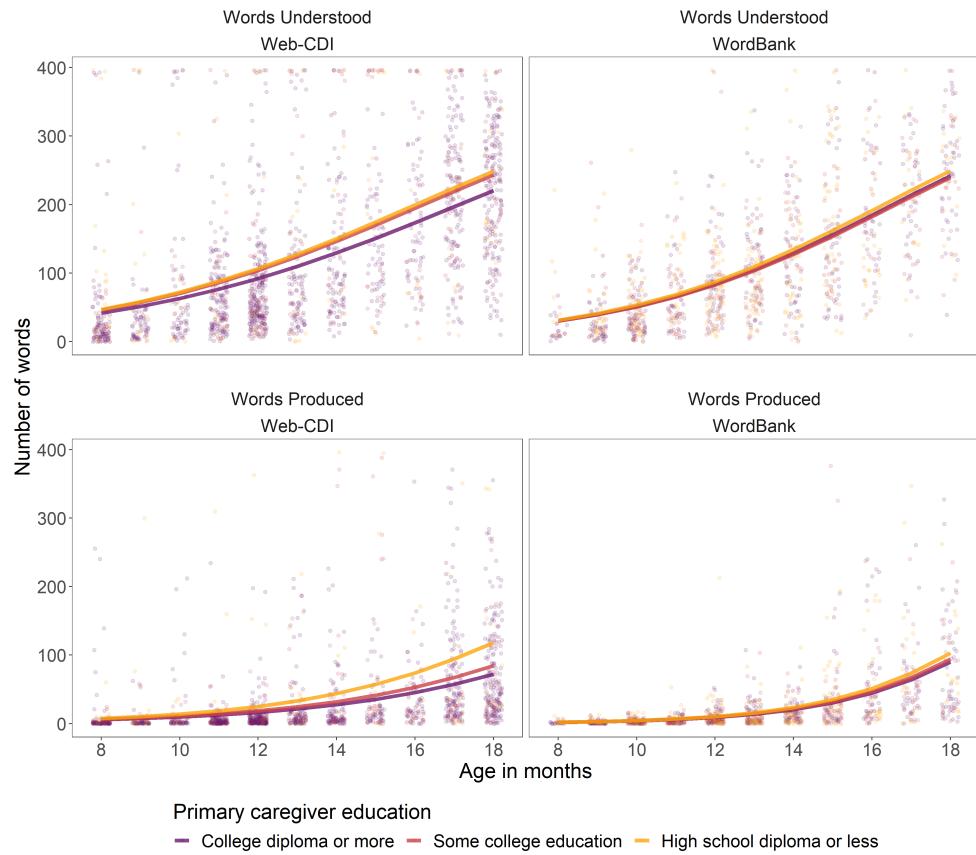


Figure 6. Individual children's word production (top panels) and comprehension (bottom panels) scores from Dataset 1 (full Web-CDI sample) plotted by age and primary caregiver's level of education (binned into "High school diploma or less," "Some college education," and "College diploma or more"). Left panels show results from the sample of Words and Gestures Web-CDI administrations collected as of November 2020 ($N = 1,620$), and right panels show the subset of American English administrations from Wordbank (Frank et al., 2021) that contain information about caregiver education ($N = 1,068$) for comparison. Curves show generalized linear model fits.

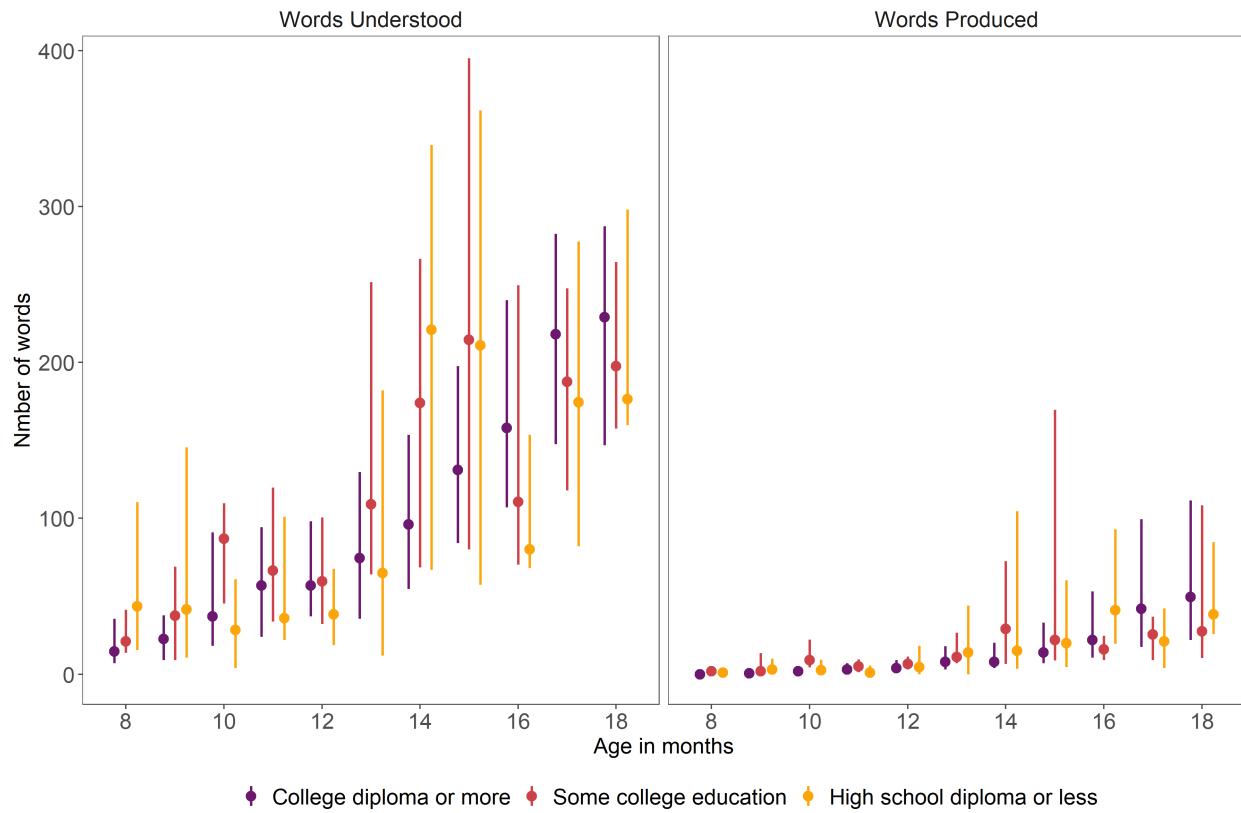


Figure 7. Median vocabulary comprehension (left) and production (right) scores from Dataset 1 (full Web-CDI sample) by age and primary caregiver's level of education attainment on the WG form. Lines indicate span between first and third quartiles for each age.

448 production scores by children's age and primary caregivers' education level shows main
 449 effects of age ($\beta = 7.60, p < 0.001$) and caregiver primary education ($\beta_{highschool} = 20.46, p$
 450 $= 0.008$). These analyses were not preregistered, but generally follow the analytic strategy
 451 in Frank, Braginsky, Yurovsky, and Marchman (2021); additionally, we fit linear models
 452 with robust standard errors to account for heteroskedasticity in the data (Astivia &
 453 Zumbo, 2019). Generalized linear model predictions for Web-CDI shown in Figure 6 differ
 454 somewhat from those for WordBank; prediction curves for caregivers of different education
 455 attainment levels diverge slightly more in the Web-CDI sample than in the WordBank
 456 sample.

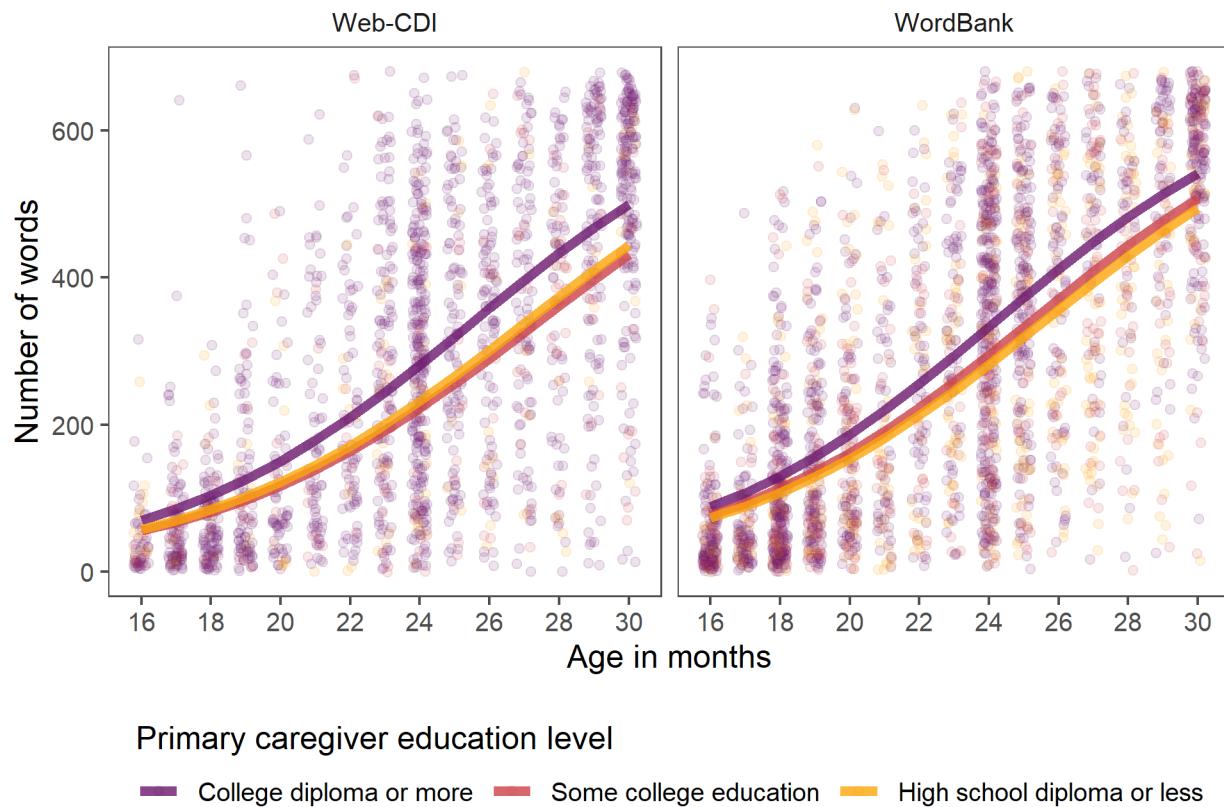


Figure 8. Individual children's vocabulary production scores from Dataset 1 (full Web-CDI sample) plotted by children's age and primary caregiver education level of primary caregiver education on as reported in the sample of Words and Sentences Web-CDI administrations collected as of November 2020 ($N = 1,900$, left panel) and in the WordBank repository ($N = 2,776$, right panel). Curves show generalized linear model fits.

457 The pattern of results seen in the WG subsample of Dataset 1 is consistent with prior
 458 findings indicating that respondents with lower levels of education attainment report
 459 higher vocabulary comprehension and production on the CDI-WG form (Feldman et al.,
 460 2000; Fenson et al., 1994). However, although caregivers with lower levels of education
 461 attainment report higher mean levels of vocabulary production and comprehension, median
 462 vocabulary scores (which are more robust to outliers) show no clear pattern of difference
 463 across primary caregiver education levels (Figure 7). This discrepancy between the

⁴⁶⁴ regression effects and a group-median analysis suggests that the regression effects described
⁴⁶⁵ previously are driven in part by differential interpretation of the survey items, such that a
⁴⁶⁶ few caregivers with lower levels of education attainment are more liberal in reporting their
⁴⁶⁷ children's production and comprehension vocabulary scores, especially for the youngest
⁴⁶⁸ children, driving up the mean scores for this demographic group.

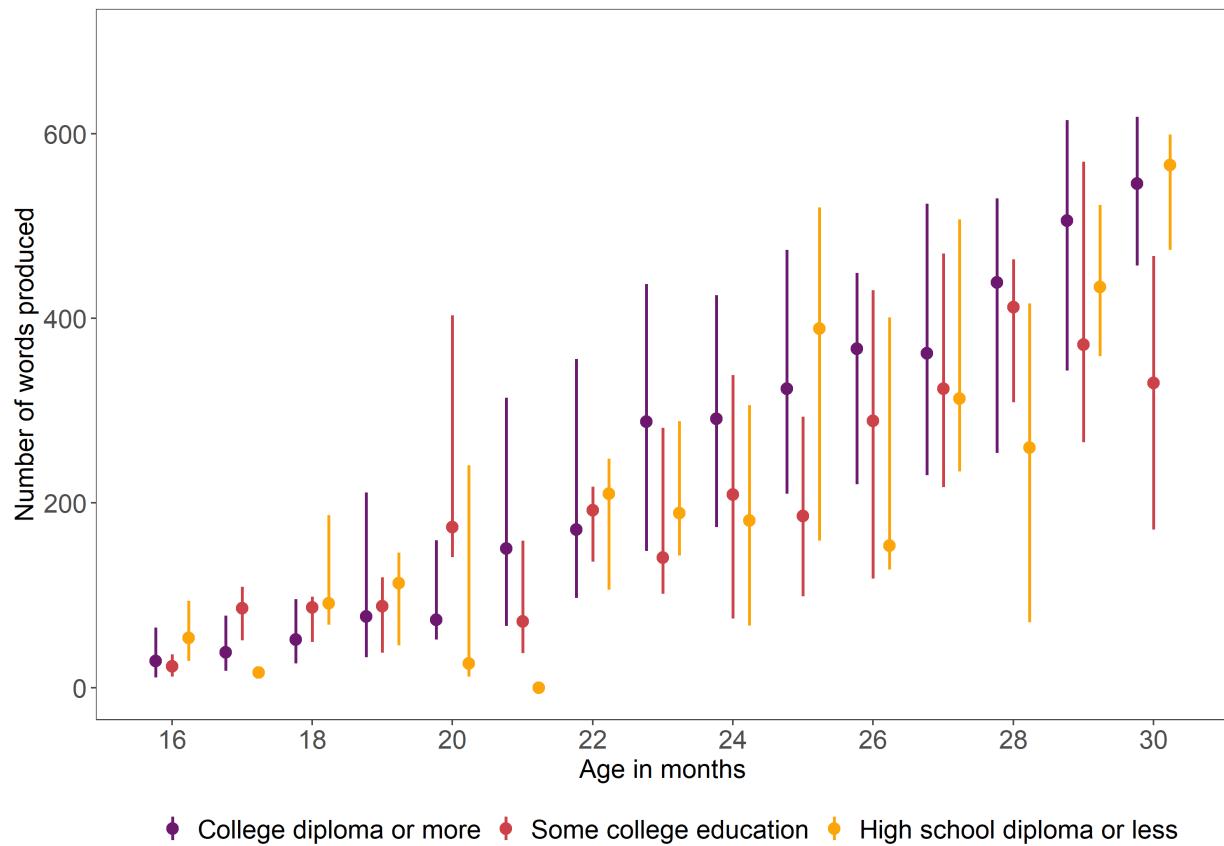


Figure 9. Median vocabulary production scores from Dataset 1 (full Web-CDI sample) by age and primary caregiver's level of education attainment on the WS form. Lines indicate span between first and third quartiles for each age.

⁴⁶⁹ Vocabulary production scores on the WS form show the expected pattern of increase
⁴⁷⁰ with children's age in months; in addition, scores replicate the trend reported in Feldman
⁴⁷¹ et al. (2000) and Frank, Braginsky, Yurovsky, and Marchman (2021) such that primary
⁴⁷² caregiver education is positively associated with children's reported vocabulary size (Figure

473 8). Because representation of caregivers without a high school diploma is scarce ($N = 6$ out
474 of a sample of 1,900), interpretation of the data from this group is constrained.

475 Nevertheless, as shown in Figure 8, a small but clear positive association between primary
476 caregiver education and vocabulary score exists such that college-educated caregivers
477 report higher vocabulary scores than those of any other education level. Notably, this
478 association is not the result of outliers and is still appreciable in median scores (Figure 9),
479 unlike the data from the WG measure shown in Figure 7. The implications from these data
480 converge with previous findings which indicate that parental education levels, often used as
481 a metric of a family's socioeconomic status, are related to children's vocabulary size
482 through early childhood.

483 **Discussion: Dataset 1.** In general, the full sample of Web-CDI data after
484 exclusions (Dataset 1) replicates previous norming datasets used with the standard
485 paper-and-pencil form of the MB-CDI. We find that vocabulary scores grow with age and
486 that females hold a slight advantage over males in early vocabulary development.
487 Moreover, Dataset 1 replicates a previously documented relationship between primary
488 caregiver education level and vocabulary scores: on the WG form, primary caregiver
489 education shows a slight negative association with vocabulary scores, whereas the trend is
490 reversed in the WS form. Taken together, these data illustrate that Web-CDI and the
491 standard paper-and-pencil form of the CDI give similar results, and thus that Web-CDI
492 can be used as a valid alternative to the paper format.

493 The data discussed above have resulted from efforts by many researchers across the
494 United States whose motivations for using the Web-CDI vary. As a result, they reproduce
495 many of the biases of standard U.S. convenience samples. In the next section, we describe
496 in more detail our recent efforts to use the Web-CDI to collect vocabulary development
497 data from traditionally underrepresented participant populations in the United States,
498 attempting to counteract these trends.

499 **Dataset 2: Using Web-CDI to Collect Data from Diverse U.S.-based
500 Communities**

501 Despite the large sample sizes we achieved in the previous section, Dataset 1 is, if
502 anything, even more biased towards highly-educated and white families than previous
503 datasets collected using the paper-and-pencil form. How can we recruit more diverse
504 samples to remedy this issue? Here, we discuss and analyze Dataset 2, which consists of
505 those administrations from Dataset 1 which were part of recent data-collection efforts
506 (within the past year and a half) that were specifically aimed towards exploring the use of
507 online recruitment as a potential way to collect more diverse participant samples than are
508 typical in the literature. In other words, the following data from Dataset 2 were included in
509 the previous discussion and analysis of Dataset 1, but we examine them separately here to
510 give special attention to the issue of collecting diverse samples online.

511 While understanding that the performance of standard measurement tools like the
512 CDI among multilinguals is of immense import to the field of vocabulary development
513 research [Gonzalez et al., in prep; Floccia et al. (2018); De Houwer (2019)], we focused in
514 Dataset 2 only on vocabulary development in monolingual children, because collecting data
515 from multilingual populations introduces additional methodological considerations (e.g.,
516 how to measure exposures in each language) that are not the focus of our work here.
517 However, it will be imperative in future to collect large-scale datasets of vocabulary data in
518 bilingual children, both to better calibrate standard tools such as the CDI, as well as to
519 reduce the bias towards monolingual families in the existing literature on measuring
520 vocabulary development.

521 **Online data collection.** Online recruitment methods, such as finding participants
522 on platforms such as Amazon Mechanical Turk, Facebook and Prolific, represent one
523 possible route towards assembling a large, diverse sample to take the Web-CDI. These
524 methods allow researchers to depart from their typical geographical recruitment area much

more easily than with paper-and-pencil administration. Online recruitment strategies for vocabulary development data collection have been used in the United Kingdom (Alcock, Meints, & Rowland, 2020), but their usage in the U.S. context remains, to our knowledge, rare. In a series of data collection efforts, we used Web-CDI as a tool to explore these different channels of recruitment.



Figure 10. Example Facebook advertisement in Phase 1 of recent data collection.

Dataset 2 consists of data that were collected in two phases. In the first phase, we ran advertisements on Facebook which were aimed at non-white families based on users' geographic locations (e.g., targeting users living in majority-Black cities) or other profile features (e.g., ethnic identification, interest in parenthood-related topics). Advertisements consisted of an image of a child and a caption informing Facebook users of an opportunity to fill out a survey on their child's language development and receive an Amazon gift card

Table 2

Exclusions from Dataset 2: recent data collection using Facebook and Prolific.

Exclusion	WG	% of full	WS	% of full
	exclusions	WG sample	exclusions	WS sample
		excluded		excluded
Not first administration	0	0.00%	0	0.00%
Premature or low birthweight	7	2.53%	1	0.33%
Multilingual exposure	18	6.50%	23	7.62%
Illnesses/Vision/Hearing	4	1.44%	4	1.32%
Out of age range	1	0.36%	26	8.61%
Completed survey too quickly	119	42.96%	133	44.04%
System error in word tabulation	0	0.00%	0	0.00%
Total exclusions	149	54%	187	62%

536 (Figure 10). Upon clicking the advertisement, participants were redirected to a unique
 537 administration of the Web-CDI; they received \$5 upon completing the survey. This
 538 open-ended approach to recruitment offered several advantages, namely that a wide variety
 539 of potential participants from specific demographic backgrounds can be reached on
 540 Facebook. However, we also received many incomplete or otherwise unusable survey
 541 administrations, either from Facebook users who clicked the link and decide not to
 542 participate, or those who completed the survey in an extremely short period of time (over
 543 half of all completed administrations, Table 2).

544 In the second phase, we used the crowdsourcing survey vendor Prolific
 545 (<http://prolific.co>) in the hopes that some of the challenges encountered with Facebook
 546 recruitment would be addressed. Prolific allows researchers to create studies and post them
 547 to individuals who are in the platform's participant database, each of whom is assigned a
 548 unique alphanumeric "Prolific ID." Importantly, Prolific maintains detailed demographic
 549 information about participants, allowing researchers to specify who they would like to

550 complete their studies. Prolific further has a built-in compensation infrastructure that
551 handles monetary payments to participants, eliminating the need to disburse gift cards
552 through Web-CDI.

553 In the particular case of Web-CDI, the demographic information needed to determine
554 whether an individual was eligible to complete our survey (e.g., has a child in the correct
555 age range, lives in a monolingual household, etc.) was more specific than the information
556 that Prolific collects about their participant base. We therefore used a brief pre-screening
557 questionnaire to generate a list of participants who were eligible to participate, and
558 subsequently advertised the Web-CDI survey to those participants. Given that we were
559 interested only in reaching participants in the United States who were not white or who
560 did not have a college diploma, our data collection efforts only yielded a sample that was
561 small ($N = 68$) but much more thoroughly screened than that which we could obtain on
562 Facebook.

563 Across both phases (Facebook and Prolific recruitment), we used the same exclusion
564 criteria as in the full Web-CDI sample to screen participants. A complete tally of all
565 excluded participants is shown in Table 2. In both the WG and WS surveys, exclusion
566 rates in Dataset 2 were high, amounting to 58% of participants who completed the survey.
567 The high exclusion rates were notably driven by an accumulation of survey administrations
568 which participants completed more quickly than our time cutoffs allow (Tables A4 and
569 A5). Many of the survey administrations excluded for fast completion also had missing
570 demographic information reported: Among WG participants excluded for too-fast
571 completions, 93% did not report ethnicity, and among WS participants excluded for the
572 same reason, 97% did not report ethnicity. Absence of these data prevents us from drawing
573 conclusions about the origin or demographic profile of administrations that were excluded.
574 After exclusions, full sample size in Dataset 2 was $N = 128$ WG completions and $N = 115$
575 WS completions.

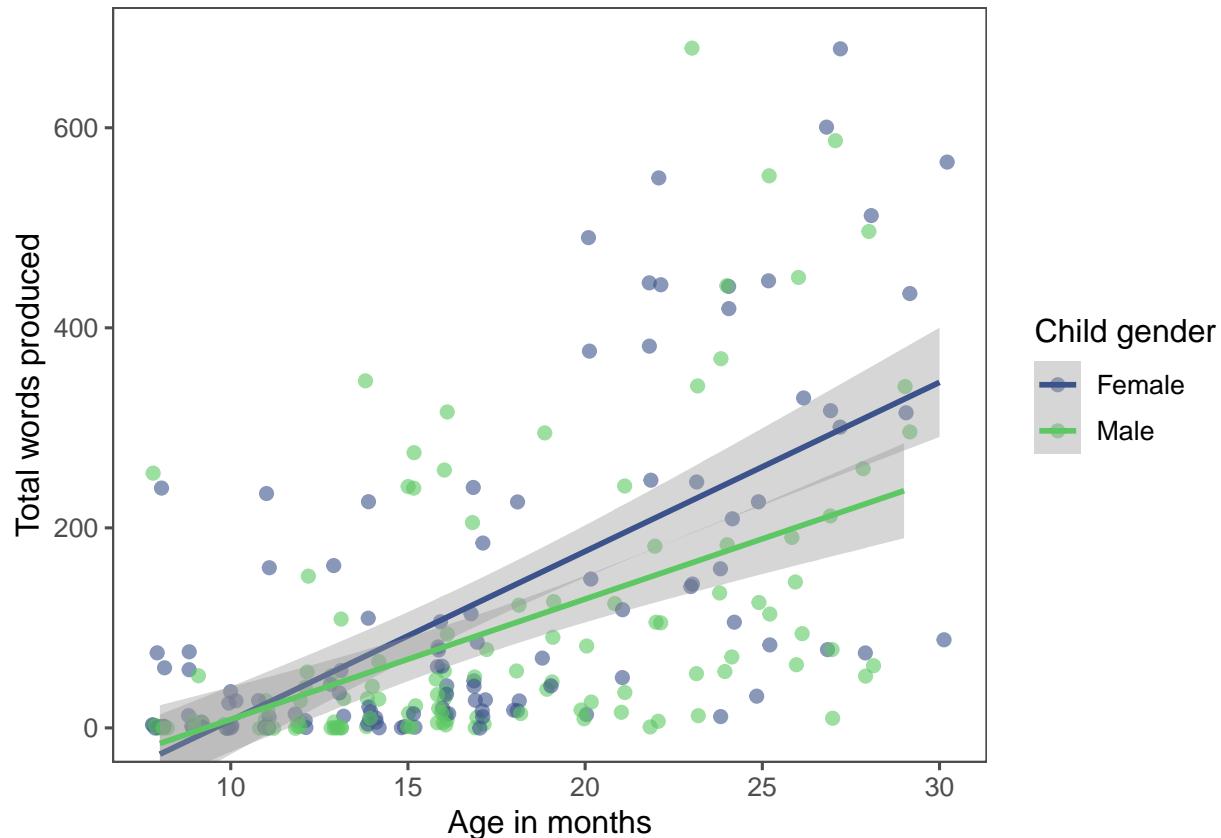


Figure 11. Individual children's vocabulary production scores from Dataset 2 (recent data collection efforts) plotted by children's age and gender (both WG and WS, N = 240, with 114 girls). Lines are best linear fits with associated 95% confidence intervals. Children with a different or no reported gender (N = 3) are omitted here.

576 The results from Dataset 2 show overall similar patterns to the full Web-CDI sample
 577 in several regards. Word production scores from both the WG and WS administrations
 578 reflect growing productive vocabulary across the second and third years, with a very small
 579 gender effect such that female children's vocabularies are higher across age than males'
 580 (Figure 11). The relationship between caregivers' reported levels of education and child's
 581 vocabulary score is not as clear as it is in the full Web-CDI sample (Figure 12); however,
 582 children of college-educated caregivers reported generally higher vocabulary scores across
 583 age than did children of caregivers without any college degree. These patterns suggest that

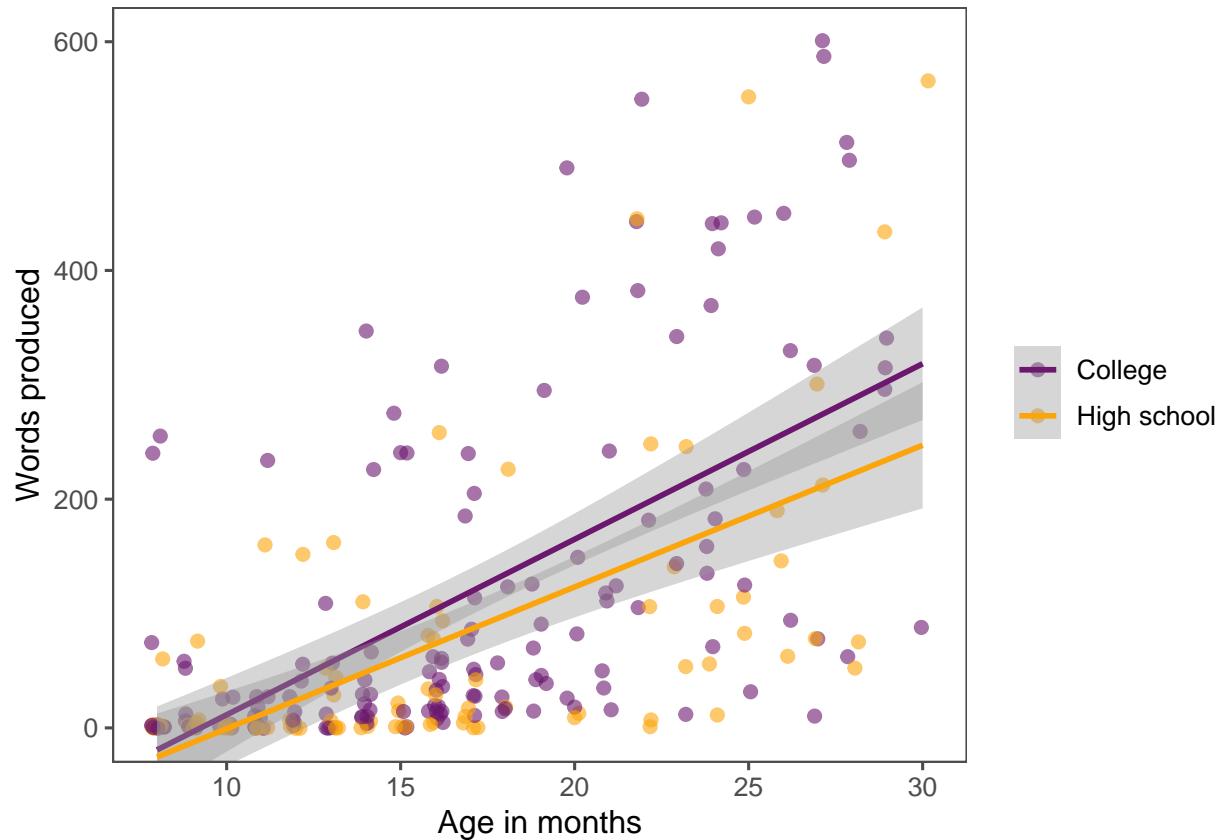


Figure 12. Individual children's vocabulary production scores from Dataset 2 (recent data collection efforts) plotted by age and level of primary caregiver education, binned into those with a high school diploma or less education and those with some college education or a college diploma ($N = 243$). Lines show best linear fits and associated 95% confidence intervals.

584 our data show similar general patterns to other CDI datasets with other populations
 585 (Frank, Braginsky, Yurovsky, & Marchman, 2021).

586 Importantly, Dataset 2 showed a substantial improvement in reaching non-white or
 587 less highly-educated participants. After exclusions, Dataset 2 has a higher proportion of
 588 non-white participants than Dataset 1 (the overall Web-CDI sample) and the norms
 589 established by Fenson et al. (2007) (Figure 13). Black participants in particular showed a
 590 marked increase in representation, from 10.5% in the 2007 norms to 30.7% in Dataset 2,

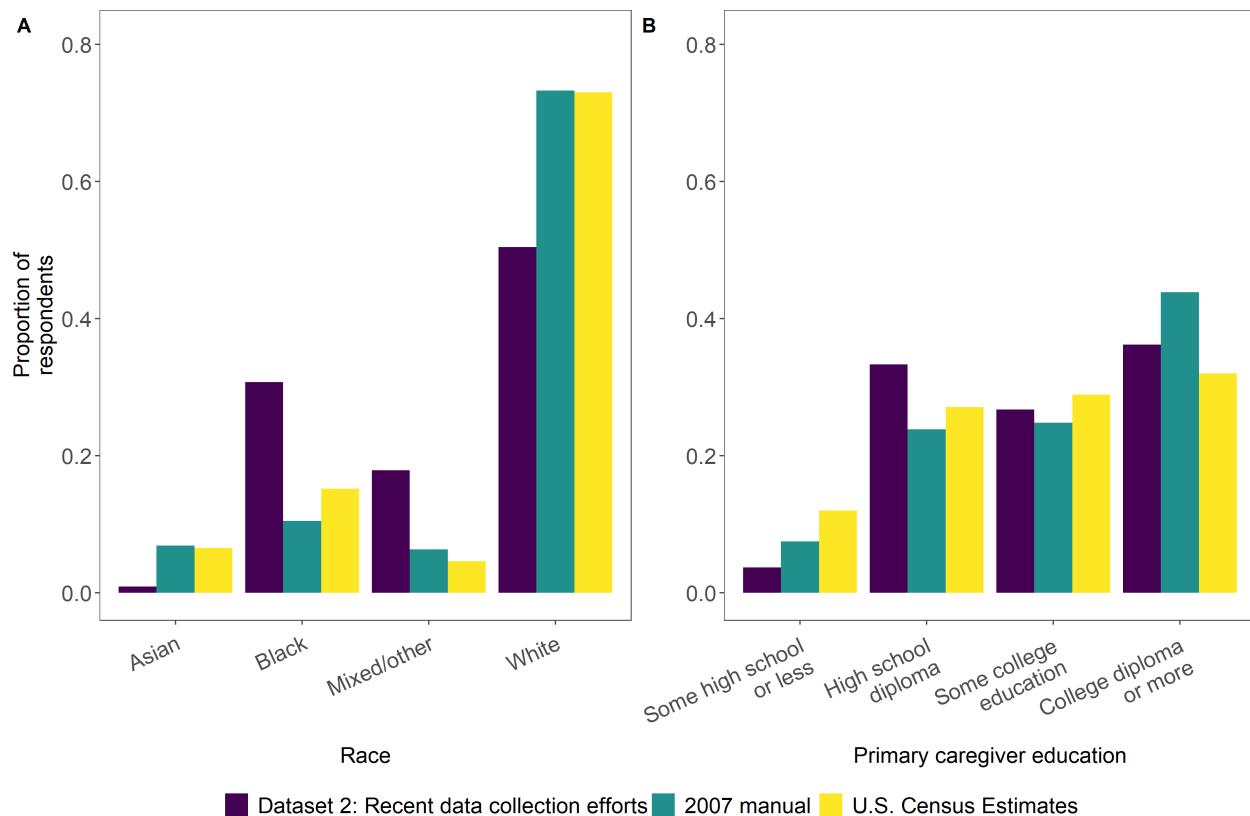


Figure 13. Proportion of respondents plotted by child race (A) and educational level of primary caregiver (B) from Dataset 2, recent data collection efforts aimed towards oversampling non-white, less highly-educated families ($N = 243$), compared with norming sample demographics from Fenson (2007). Latinx participants can be of any race and are thus not represented as a separate category here.

591 while the proportion of white participants decreased from 73.3% in the 2007 norms to
 592 50.5% in Dataset 2. Representation on the basis of families' reported primary caregiver
 593 education also improved (Figure 13). Participants with only a high school diploma
 594 accounted for 33.3% of Dataset 2 as compared to 23.8% in the 2007 norms, and
 595 representation of those with a college diploma or more education decreased from 43.8% in
 596 the 2007 norms to 36.2% in Dataset 2. Notably, the distribution of Dataset 2 with regards
 597 to primary caregiver education level is quite similar to Kristoffersen et al. (2013), who

598 collected a large, nationally-representative sample of CDI responses in Norway and
599 obtained a sample with 30%, 42%, and 24% for participants reporting 12, 14-16, and 16+
600 years of education, respectively.

601 **Discussion: Dataset 2.** The results from Dataset 2 indicate that Web-CDI could
602 be a promising platform to collect vocabulary development data in non-white populations
603 and communities with lower levels of education attainment when paired with online
604 recruitment methods that yield legitimate, representative participant samples. At the same
605 time, however, these data convey clear limitations of our approach. Perhaps most
606 conspicuously, more than half of completed administrations in this sample had to be
607 excluded, in many cases because the information provided by participants appeared rushed
608 or incomplete: over 40% of administrations were completed in a shorter amount of time
609 than that allowed by our cutoff criteria (Tables A4 and A5), and of these quick
610 completions, well over 90% were missing demographic information that is rarely missing in
611 other administrations of the form. Determining the precise reasons for the high exclusion
612 rate, and how (if at all) this (self-)selection may bias data reflecting demographic trends in
613 vocabulary development, requires a more thorough assessment of who is submitting
614 hastily-completed forms. Such an assessment is beyond the scope of the current study.
615 However, all respondents who got to the end of the form were compensated regardless of
616 how thoroughly they completed it, creating the possibility that some participants who
617 clicked the anonymous link may not have been members of the population of interest, but
618 rather were other individuals motivated by compensation. To the extent that participants
619 moved through the form quickly because they found the length burdensome, a transition to
620 short forms, including computer adaptive ones (e.g., Chai, Lo, & Mayor, 2020; Kachergis et
621 al., 2021; Makransky, Dale, Havmose, & Bleses, 2016; Mayor & Mani, 2019), would
622 potentially increase data quality and completion rates substantially.

623 Additionally, the exclusion rates described previously provide information only on
624 those participants who did, at some point, submit a completed form, but many individuals

625 clicked the advertisement link and did not subsequently continue on to complete the form.
626 Without an in-depth exploration of who is clicking the link and why they might choose not
627 to continue, we cannot draw conclusions about the representativeness of the sample in
628 Dataset 2 with regards to the communities we would like to include in our research. As
629 such, a more thorough understanding of how users from different communities respond to
630 various recruitment and sampling methods is needed in future work in order to draw
631 conclusions about demographic trends above and beyond those already established in the
632 literature.

633 Participants in Dataset 2 were recruited through a targeted post on social media, a
634 technique that is considerably more anonymous than recruitment strategies which entail
635 face-to-face or extended contact between researchers and community members. Online
636 recruitment methods may not be suitable for all communities, especially when researchers
637 ask participants to report potentially sensitive information about the health, developmental
638 progress, ethnicity and geographic location of their children (even when such information is
639 stored anonymously). Our goal here was to assess whether general trends in past literature
640 could be recovered using such an online strategy, but future research should take into
641 account that other more personal methods of recruitment, such as direct community
642 outreach or liaison contacts, may improve participants' experiences and their willingness to
643 engage with the study.

644 Finally, a significant limitation of the data collection process in Dataset 2 is that
645 many people in the population of interest - particularly lower-income families - do not have
646 reliable internet access. Having participants complete the Web-CDI on a mobile device
647 may alleviate some of the issues caused by differential access to Wi-Fi, since the vast
648 majority of American adults own a smartphone (Pew Research Center, 2019). Accordingly,
649 improving Web-CDI's user experience on mobile platforms will be an important step
650 towards ensuring that caregivers across the socioeconomic spectrum can easily complete
651 the survey. For smartphone users on pay-as-you-go plans, who may be reluctant to use

652 phone data to complete a study, a possible solution could be compensating participants for
653 the amount of “internet time” they incurred completing the form.

654 **General Discussion and Conclusions**

655 In this paper, we have presented Web-CDI, a comprehensive online interface for
656 researchers to measure children’s vocabulary by administering the MacArthur-Bates
657 Communicative Development Inventories family of parent-report instruments. Web-CDI
658 provides a convenient researcher management interface, built-in data privacy protections,
659 and a variety of features designed to make both longitudinal and social-media sampling
660 easy. To date, over 3,500 valid administrations of the WG and WS forms have been
661 collected on Web-CDI from more than a dozen researchers in the United States after
662 applying strict exclusion criteria derived from previous norming studies (Fenson et al., 2007,
663 1994). Our analysis of Dataset 1 shows that demographic trends from previous work using
664 the paper-and-pencil CDI form are replicated in data gleaned from Web-CDI, suggesting
665 that the Web-CDI is a valid alternative to the paper form and captures similar results.

666 Many research laboratories, not only in the United States but around the world,
667 collect vocabulary development data using the MacArthur-Bates CDI in its original or
668 adapted form. With traditional paper-based forms, combining insights from various
669 research groups can prove challenging, as each group may have slightly different ways of
670 formatting and managing data from CDI forms. By contrast, if all of these groups’ data
671 come to be stored in a single repository with a consistent database structure, data from
672 disparate sources can easily be collated and analyzed in a uniform fashion. As such, a
673 centralized repository such as Web-CDI provides a streamlined data-aggregation pipeline
674 that facilitates cross-lab collaborations, multisite research projects and the curation of large
675 datasets that provide more power to characterize the vast individual differences present in
676 children’s vocabulary development.

677 Beyond the goal of simply getting more data, we hope that Web-CDI can advance
678 efforts to expand the reach of vocabulary research past convenience samples into diverse
679 communities. A key question in the field of vocabulary development concerns the
680 mechanisms through which sociodemographic variables, such as race, ethnicity, income and
681 education are linked to group differences in vocabulary outcomes. Large,
682 population-representative samples of vocabulary development data are needed to
683 understand these mechanisms, but research to date (including the full sample of Web-CDI
684 administrations) has often oversampled non-Hispanic white participants and those with
685 advanced levels of education.

686 We explored the use of Web-CDI as part of a potential strategy to collect data from
687 non-white and less highly-educated communities in two phases (Dataset 2). Several overall
688 patterns emerged which we expected: vocabulary scores grew with age, providing a basic
689 validity check of the Web-CDI measure; females held a slight advantage in word learning
690 over males; and children of caregivers with a college education showed slightly higher
691 vocabulary scores. Nonetheless, the insights from these data, while aligned with past
692 norming studies, are necessarily constrained by several features of our method.

693 Limitations of our method notwithstanding, a transition to web-based data collection
694 streamlines the process by which historically underrepresented populations can be reached
695 in child language research. In particular, recruitment methods involving community
696 partners, such as parenting groups, childcare centers and early education providers, are
697 simplified substantially if leaders in these organizations can distribute a web survey to their
698 members that is easy to fill out, as compared with paper forms, which typically present
699 logistical hurdles for distribution and collection. Additionally, we hope that Web-CDI can
700 serve as an accessible, free, and easy to use resource for researchers already doing extensive
701 work with underrepresented groups.

702 Web-based data collection can capture useful information about vocabulary

703 development from diverse communities, but future research will need to examine which
704 sampling methods can yield accurate, population-representative data that can advance our
705 understanding of the link between sociodemographic variation and variation in language
706 outcomes.

707 **Acknowledgements**

708 We thank Larry Fenson, Philip Dale, and Donna Thal for their assistance and helpful
709 feedback preparing this manuscript.

710 **Ethics statement**

711 Data collected in the United States for this project are anonymized according to
712 guidelines set forth by the United States Department of Health and Human Services. Data
713 collection at Stanford University was approved by the Stanford Institutional Review Board
714 (IRB), protocol 20398.

715 **Data, code and materials availability statement**

- 716 • Open data: All data analyzed in this work are available on the Open Science
717 Framework at <https://osf.io/nmdq4/>.
- 718 • Code: All code for this work is available on the Open Science Framework at
719 <https://osf.io/nmdq4/>.
- 720 • Materials: All code and materials for the Web-CDI are openly available at
721 <https://github.com/langcog/web-cdi>. If readers wish to view the Web-CDI interface
722 in full from the participants' or researchers' perspectives, they are encouraged to
723 contact webcdi-contact@stanford.edu.

Author contributions

- Conceptualization: Benjamin deMayo, Danielle Kellier, Mika Braginsky, Christina Bergmann, Caroline Rowland, Michael Frank and Virginia Marchman.
- Data Curation: Benjamin deMayo, Danielle Kellier and Virginia Marchman.
- Formal Analysis: Benjamin deMayo.
- Funding Acquisition: Caroline Rowland and Michael Frank.
- Investigation: Benjamin deMayo, Danielle Kellier and Virginia Marchman.
- Methodology: Benjamin deMayo, Danielle Kellier, Michael Frank and Virginia Marchman.
- Project Administration: Caroline Rowland, Michael Frank and Virginia Marchman.
- Software: Danielle Kellier, Mika Braginsky, Christina Bergmann and Cielke Hendriks.
- Supervision: Christina Bergmann, Caroline Rowland, Michael Frank and Virginia Marchman.
- Visualization: Benjamin deMayo.
- Writing - Original Draft Preparation: Benjamin deMayo, Michael Frank and Virginia Marchman.
- Writing - Review & Editing: Benjamin deMayo, Danielle Kellier, Mika Braginsky, Christina Bergmann, Cielke Hendriks, Caroline Rowland, Michael Frank and Virginia Marchman.

Software used

R [Version 4.0.3; R Core Team (2020)] and the R-packages *broman* [Version 0.71.6; Broman (2020)], *cowplot* [Version 1.1.0; Wilke (2020)], *dplyr* [Version 1.0.2; Wickham, François, Henry, and Müller (2020)], *estimatr* [Version 0.26.0; Blair, Cooper, Coppock, Humphreys, and Sonnet (2020)], *forcats* [Version 0.5.0; Wickham (2020a)], *fs* [Version 1.5.0; Hester and Wickham (2020)], *ggplot2* [Version 3.3.2; Wickham (2016)], *here* [Version 0.1; Müller (2017)], *kableExtra* [Version 1.3.1; Zhu (2020)], *papaja* [Version 0.1.0.9997; Aust and

750 Barth (2020)], *purrr* [Version 0.3.4; Henry and Wickham (2020)], *readr* [Version 1.4.0;
751 Wickham and Hester (2020)], *scales* [Version 1.1.1; Wickham and Seidel (2020)], *stringr*
752 [Version 1.4.0; Wickham (2019)], *tibble* [Version 3.0.4; Müller and Wickham (2020)], *tidyR*
753 [Version 1.1.2; Wickham (2020b)], *tidyverse* [Version 1.3.0; Wickham et al. (2019)],
754 *wordbankr* [Version 0.3.1; (R-wordbankr?)], and *xtable* [Version 1.8.4; Dahl, Scott,
755 Roosen, Magnusson, and Swinton (2019)]

References

- 756 Alcock, K., Meints, K., & Rowland, C. (2020). *The UK communicative development inventories: Words and gestures*. J&R Press.
- 757
- 758
- 759 Astivia, O. L. O., & Zumbo, B. D. (2019). Heteroskedasticity in multiple regression analysis: What it is, how to detect it and how to solve it with applications in r
- 760
- 761 and SPSS. *Practical Assessment, Research, and Evaluation*, 24(1), 1.
- 762 Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*.
- 763 *Journal of Open Source Software* (Vol. 4, p. 1686).
- 764 <https://doi.org/10.21105/joss.01686>
- 765 Bates, E., & Goodman, J. C. (2001). On the inseparability of grammar and the
- 766 lexicon: Evidence from acquisition.
- 767 Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P., Reznick, J. S., ... Hartung,
- 768 J. (1994). Developmental and stylistic variation in the composition of early
- 769 vocabulary. *J Child Lang*, 21(01), 85–123.
- 770 Blair, G., Cooper, J., Coppock, A., Humphreys, M., & Sonnet, L. (2020). *Estimatr: Fast estimators for design-based inference*. *Journal of Open Source Software* (Vol.
- 771 4, p. 1686). Springer-Verlag New York. <https://doi.org/10.21105/joss.01686>
- 772
- 773 Bornstein, M. H., & Putnick, D. L. (2012). Stability of language in childhood: A
- 774 multiage, multidomain, multimeasure, and multisource study. *Developmental*
- 775 *Psychology*, 48(2), 477.
- 776 Broman, K. W. (2020). *Broman: Karl broman's r code*. *Journal of Open Source*
- 777 *Software* (Vol. 4, p. 1686). Springer-Verlag New York.
- 778 <https://doi.org/10.21105/joss.01686>
- 779 Caselli, N. K., Lieberman, A. M., & Pyers, J. E. (2020). The ASL-CDI 2.0: An
- 780 updated, normed adaptation of the MacArthur bates communicative

- development inventory for american sign language. *Behavior Research Methods*, 1–14.
- Chai, J. H., Lo, C. H., & Mayor, J. (2020). A bayesian-inspired item response theory-based framework to produce very short versions of MacArthur–bates communicative development inventories. *Journal of Speech, Language, and Hearing Research*, 63(10), 3488–3500.
- Dahl, D. B., Scott, D., Roosen, C., Magnusson, A., & Swinton, J. (2019). *Xtable: Export tables to LaTeX or HTML*. Retrieved from <https://CRAN.R-project.org/package=xtable>
- Dale, P. S. (2015). Adaptations, Not Translations! Retrieved from <http://mb-cdi.stanford.edu/Translations2015.pdf>
- De Houwer, A. (2019). Equitable evaluation of bilingual children’s language knowledge using the CDI: It really matters who you ask. *Journal of Monolingual and Bilingual Speech*, 1(1), 32–54.
- Feldman, H. M., Dollaghan, C. A., Campbell, T. F., Kurs-Lasky, M., Janosky, J. E., & Paradise, J. L. (2000). Measurement properties of the MacArthur communicative development inventories at ages one and two years. *Child Development*, 71(2), 310–322.
- Fenson, L., Bates, E., Dale, P. S., Marchman, V. A., Reznick, J. S., & Thal, D. J. (2007). *MacArthur-Bates Communicative Development Inventories*. Brookes Publishing Company.
- Fenson, L., Dale, P. S., Reznick, J. S., Bates, E., Thal, D. J., Pethick, S. J., ... Stiles, J. (1994). Variability in early communicative development. *Monogr Soc Res Child Dev*, 59(5).

- 805 Fenson, L., Pethick, S., Renda, C., Cox, J. L., Dale, P. S., & Reznick, J. S. (2000).
806 Short-form versions of the MacArthur communicative development inventories.
807 *Applied Psycholinguistics*, 21(1), 95–116.
- 808 Floccia, C., Sambrook, T. D., Delle Luche, C., Kwok, R., Goslin, J., White, L., ...
809 others. (2018). Vocabulary of 2-year-olds learning english and an additional
810 language: Norms and effects of linguistic distance. *Monographs of the Society for*
811 *Research in Child Development*, 83(1), 1–135.
- 812 Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., ...
813 others. (2017). A collaborative approach to infant research: Promoting
814 reproducibility, best practices, and theory-building. *Infancy*, 22(4), 421–435.
- 815 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability*
816 *and consistency in early language learning: The wordbank project*. MIT Press.
- 817 Henry, L., & Wickham, H. (2020). *Purrr: Functional programming tools*. *Journal of*
818 *Open Source Software* (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>
- 819 Hester, J., & Wickham, H. (2020). *Fs: Cross-platform file system operations based*
820 *on 'libuv'*. *Journal of Open Source Software* (Vol. 4, p. 1686). Springer-Verlag
821 New York. <https://doi.org/10.21105/joss.01686>
- 822 Kachergis, G., Marchman, V., Dale, P., Mehta, H., Mankewitz, J., & Frank, M.
823 (2021). *An online computerized adaptive test (CAT) of children's vocabulary*
824 *development in english and mexican spanish*.
- 825 Kapalková, S., & Slanèová, D. (2007). Adaptation of CDI to the slovak language. In
826 *Proceedings from the first european network meeting on the communicative*
827 *development inventories: May 24-28 2006 dubrovnik croatia*. University of Gävle.
- 828 Kartushina, N., Mani, N., AKTAN-ERCIYES, A., Alaslani, K., Aldrich, N. J.,
829 Almohammadi, A., ... al., et. (2021). COVID-19 first lockdown as a unique

- 830 window into language acquisition: What you do (with your child) matters.
- 831 PsyArXiv. <https://doi.org/10.31234/osf.io/5ejwu>
- 832 Kristoffersen, K. E., Simonsen, H. G., Bleses, D., Wehberg, S., Jørgensen, R. N.,
- 833 Eiesland, E. A., & Henriksen, L. Y. (2013). The use of the internet in collecting
- 834 CDI data—an example from norway. *Journal of Child Language*, 40(03), 567–585.
- 835 Makransky, G., Dale, P. S., Havmose, P., & Bleses, D. (2016). An item response
- 836 theory-based, computerized adaptive testing version of the MacArthur–bates
- 837 communicative development inventory: Words & sentences (CDI: WS). *Journal*
- 838 *of Speech, Language, and Hearing Research*, 59(2), 281–289.
- 839 Mayor, J., & Mani, N. (2019). A short version of the MacArthur–bates
- 840 communicative development inventories with high validity. *Behavior Research*
- 841 *Methods*, 51(5), 2248–2255.
- 842 Müller, K. (2017). *Here: A simpler way to find your files*. *Journal of Open Source*
- 843 *Software* (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>
- 844 Müller, K., & Wickham, H. (2020). *Tibble: Simple data frames*. *Journal of Open*
- 845 *Source Software* (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>
- 846 Percentage of persons 18 to 24 years old and age 25 and over, by educational
- 847 attainment, race/ethnicity, and selected racial/ethnic subgroups: 2010 and 2017.
- 848 (2019). https://nces.ed.gov/programs/digest/d18/tables/dt18_104.40.asp?referer=raceindica.asp.
- 849
- 850 Pew research center mobile fact sheet. (2019). Retrieved from
- 851 <https://www.pewresearch.org/internet/fact-sheet/mobile/>
- 852 R Core Team. (2020). *R: A language and environment for statistical computing*.
- 853 *Journal of Open Source Software* (Vol. 4, p. 1686). Vienna, Austria: R
- 854 Foundation for Statistical Computing; Springer-Verlag New York.

- 855 <https://doi.org/10.21105/joss.01686>
- 856 Snyder, T. D., De Brey, C., & Dillow, S. A. (2019). Digest of education statistics
857 2017, NCES 2018-070. *National Center for Education Statistics*.
- 858 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. *Journal of Open*
859 *Source Software* (Vol. 4, p. 1686). Springer-Verlag New York.
860 <https://doi.org/10.21105/joss.01686>
- 861 Wickham, H. (2019). *Stringr: Simple, consistent wrappers for common string*
862 *operations*. *Journal of Open Source Software* (Vol. 4, p. 1686).
863 <https://doi.org/10.21105/joss.01686>
- 864 Wickham, H. (2020a). *Forcats: Tools for working with categorical variables*
865 *(factors)*. *Journal of Open Source Software* (Vol. 4, p. 1686). Springer-Verlag
866 New York. <https://doi.org/10.21105/joss.01686>
- 867 Wickham, H. (2020b). *Tidyr: Tidy messy data*. *Journal of Open Source Software*
868 (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>
- 869 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ...
870 Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*,
871 4(43), 1686. <https://doi.org/10.21105/joss.01686>
- 872 Wickham, H., François, R., Henry, L., & Müller, K. (2020). *Dplyr: A grammar of*
873 *data manipulation*. *Journal of Open Source Software* (Vol. 4, p. 1686).
874 Springer-Verlag New York. <https://doi.org/10.21105/joss.01686>
- 875 Wickham, H., & Hester, J. (2020). *Readr: Read rectangular text data*. *Journal of*
876 *Open Source Software* (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>
- 877 Wickham, H., & Seidel, D. (2020). *Scales: Scale functions for visualization*. *Journal*
878 *of Open Source Software* (Vol. 4, p. 1686). <https://doi.org/10.21105/joss.01686>

- 879 Wilke, C. O. (2020). *Cowplot: Streamlined plot theme and plot annotations for*
880 *'ggplot2'*. *Journal of Open Source Software* (Vol. 4, p. 1686). Springer-Verlag
881 New York. <https://doi.org/10.21105/joss.01686>
- 882 Zhu, H. (2020). *kableExtra: Construct complex table with 'kable' and pipe syntax.*
883 *Journal of Open Source Software* (Vol. 4, p. 1686).
884 <https://doi.org/10.21105/joss.01686>

Appendix

Table A1

Settings customizable by researchers when creating new studies to be run on the Web-CDI platform.

Study setting	Default value	Notes
Study name	none	—
Instrument	none	—
Age range for study	none	Defaults based on instrument selected.
Number of days before study expiration	14	Must be between 1 and 28 days.
Measurement units for birth weight	Pounds and ounces	Weight can also be measured in kilograms (kg).
Minimum time (minutes) a parent must take to complete the study	6	—
Waiver of documentation	blank	Can be filled in by researchers to include a Waiver of Documentation for the participant to approve before proceeding to the experiment.
Pre-fill data for longitudinal participants?	No, do not populate any part of the form	Researchers can choose to pre-fill the background information and the vocabulary checklist.

Table A1

Settings customizable by researchers when creating new studies to be run on the Web-CDI platform. (continued)

Study setting	Default value	Notes
Would you like to pay subjects in the form of Amazon gift cards?	No	If checked, researchers can enter gift codes to distribute to participants once they have completed the survey.
Do you plan on collecting only anonymous data in this study? (e.g., posting ads on social media, mass emails, etc)	No	If checked, researchers can set a limit for the maximum number of participants, as well as select an option that asks participants to verify that the information entered is accurate.
Would you like to show participants graphs of their data after completion?	Yes	–
Would you like participants to be able to share their Web-CDI results via Facebook?	No	–
Would you like participants to answer the confirmation questions?	No	Asks redundant demographic questions to serve as attention checks.

Table A1

Settings customizable by researchers when creating new studies to be run on the Web-CDI platform. (continued)

Study setting	Default value	Notes
Provide redirect button at completion of study?	No	Used to redirect users to external site after form completion.
Capture the Prolific Id for the participant?	No	For integration with Prolific.
Allow participant to print their responses at end of Study?	No	—
End message	Standard end-of-study message	Can be changed to customize end-of-study message.

Table A2

Regression output for WG comprehension measure.

term	estimate	std.error	statistic	p.value	conf.low	conf.high	df
Intercept	122.275	2.427	50.381	0.000	117.515	127.035	1610
Age	20.050	0.767	26.127	0.000	18.545	21.556	1610
Caregiver education: Some college	17.445	8.179	2.133	0.033	1.403	33.487	1610
Caregiver education: High school or less	21.862	10.935	1.999	0.046	0.413	43.311	1610
Age * Caregiver education: Some college	-1.991	2.261	-0.881	0.379	-6.425	2.443	1610
Age * Caregiver education: High school or less	-6.604	3.159	-2.091	0.037	-12.800	-0.408	1610

Table A3

Regression output for WG production measure.

term	estimate	std.error	statistic	p.value	conf.low	conf.high	df
Intercept	29.771	1.332	22.358	0.000	27.159	32.382	1610
Age	7.599	0.498	15.264	0.000	6.622	8.575	1610
Caregiver education: Some college	5.640	4.919	1.147	0.252	-4.009	15.289	1610
Caregiver education: High school or less	20.455	7.693	2.659	0.008	5.366	35.545	1610
Age * Caregiver education: Some college	-1.357	1.327	-1.022	0.307	-3.960	1.247	1610
Age * Caregiver education: High school or less	-0.121	2.095	-0.058	0.954	-4.229	3.988	1610

Table A4

Minimum times to completion, WG measure

Age in months	Minimum time to completion (minutes)
8	3.496
9	4.057
10	4.619
11	5.181
12	5.743
13	6.305
14	6.867
15	7.429
16	7.991
17	8.553
18	9.115

Table A5

Minimum times to completion, WG measure

Age in months	Minimum time to completion (minutes)
16	8.129
17	8.613
18	9.097
19	9.581
20	10.065
21	10.55
22	11.034
23	11.518
24	12.002
25	12.486
26	12.97
27	13.455
28	13.939
29	14.423
30	14.907