

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

Benjamin deMayo<sup>1</sup>, Danielle Kellier<sup>2</sup>, Mika Braginsky<sup>3</sup>, Christina Bergmann<sup>4</sup>, Cielke  
Hendriks<sup>4</sup>, Caroline Rowland<sup>4</sup>, Michael C. Frank<sup>5</sup>, & Virginia Marchman<sup>5</sup>

<sup>1</sup> Princeton University

<sup>2</sup> University of Pennsylvania

<sup>3</sup> Massachusetts Institute of Technology

<sup>4</sup> Max Planck Institute for Psycholinguistics

<sup>5</sup> Stanford University

## Abstract

Understanding the mechanisms that drive variation in children’s language acquisition requires large, population-representative datasets of children’s word learning across development. Parent report measures such as the MacArthur-Bates Communicative Development Inventories (CDI) are commonly used to collect such data, but the traditional paper-based forms make the curation of large datasets logistically challenging. Many CDI datasets are thus gathered using convenience samples, often recruited from communities in proximity to major research institutions. Here, we introduce Web-CDI, a web-based tool which allows researchers to collect CDI data online. Web-CDI contains functionality to collect and manage longitudinal data, share links, and download standardized vocabulary scores. To date, over 3,500 valid Web-CDI administrations have been completed. General trends found in past norming studies of the CDI are present in data collected from Web-CDI: scores of children’s productive vocabulary grow with age, female children show a slightly faster rate of vocabulary growth in early childhood, and participants with higher levels of educational attainment report slightly higher vocabulary production scores. We also report results from an effort to oversample non-white, lower-SES participants via online recruitment ( $N = 241$ ). These data showed similar demographic trends to the full sample but this effort recruited in a high exclusion rate. We conclude by discussing implications and challenges for the collection of large, population-representative datasets.

*Keywords:* vocabulary development, parent report, socioeconomic status

Word count: X

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

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

In this paper, we introduce a web-based tool, Web-CDI, which was developed to address the need for collecting CDI data in an online format. Web-CDI allows researchers to increase the convenience of CDI administration, further decrease costs associated with data collection and entry, and access participant samples that have traditionally been difficult to reach in language development research. Our purpose in this paper is twofold: first, we describe Web-CDI as a platform which streamlines the process of collecting MB-CDI data and collates the data in a way that facilitates the creation of large-scale, multisite collaborative datasets. Second, we profile usage of Web-CDI thus far, with a particular focus on broadening the reach of traditional paper-based methods of collecting vocabulary development data.

## The Importance of Parent Report Data

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

daily experiences with the child, which are much more extensive than a researcher or clinician can generally obtain. Moreover, they are less likely to be influenced by factors that may mask a child's true ability in the laboratory or clinic (e.g., shyness). One widely used set of parent-report instruments is the MacArthur-Bates Communicative Development Inventories, originally designed for children learning American English (Fenson et al., 2007). The American English CDIs come in two versions, Words & Gestures for children 8 to 18 months, focusing on word comprehension and production, as well as gesture use, and Words & Sentences, for children 16 to 30 months, focusing on word production and sentence structure. Together, these instruments allow for a comprehensive picture of milestones that characterize language development in early childhood.

A substantial body of evidence suggests that these instruments are both reliable and valid (e.g., Fenson et al., 1994, 2007) leading to their widespread use in thousands of research studies over the last few decades. Indeed, the popularity of the American English and Spanish CDI instruments has meant that many teams around the world have adapted the CDI format to the particular language and community (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 (e.g., Cheerios are more common in American than British homes, so age of acquisition of this word may differ substantially). To date there are now more than 100 adaptations for languages around the globe.

Initial large-scale work to establish the normative datasets for the American English CDI not only provided key benchmarks for determining children's progress, but also documented the extensive individual differences that characterize early language learning during this critical period of development (Bates et al., 1994; Fenson et al., 1994). Understanding the origins and consequences of this variability remains an important empirical and theoretical endeavor (e.g., Bates & Goodman, 2001; Bornstein & Putnick, 2012; see also, Frank et al., 2021). The popularity of CDI instruments has remained strong

over the years, leading to extensions of the methodology to alternative formats, e.g., short forms (Fenson et al., 2000).

While the reliability and validity of these instruments is well-established for the American English versions of the forms, existing norming samples are skewed toward families with more years of formal education and away from non-White groups (Fenson et al., 2007). Representation in these norming samples is generally restricted to families living on the US 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 MacArthur-Bates 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 using a custom online platform. Similarly, the Slovak adaptation of the CDI uses an online administration format. And many groups have used general purpose survey software such as Qualtrics and Survey Monkey to administer CDIs and variants online (e.g., Caselli, Lieberman, & Pyers, 2020). The innovation of Web-CDI is to provide a comprehensive researcher management interface for the administration of a wide range of CDI forms, allowing researchers to manage longitudinal administrations, download standardized scores, and share links easily, all while satisfying strong guarantees regarding privacy and anonymity. Moreover, a key benefit of a unified data collection and storage system such as Web-CDI is that data from disparate sources are combined into a single repository, substantially reducing the overhead efforts associated with bringing together data collected using paper forms by researchers across the world.

## Introducing Web-CDI

Web-CDI is a web-based platform for CDI administration and management. Web-CDI allows researchers to communicate with families by sharing URLs via email or social media, facilitating access to families in areas distant from an academic institution and eliminating costly mailings and laboratory visits. Web-CDI also standardizes electronic administration and scoring of CDI forms across labs and institutions, making possible the aggregation of CDI data for later reuse and comparison across administrations by different labs. Indeed, users of Web-CDI grant the CDI Advisory Board to access and analyze the resulting data on an opt-out basis, providing a path towards continual improvement of CDI instruments. Since 2018, more than 3,500 CDIs have been collected by 15 research groups throughout the US who are using Web-CDI, demonstrating the potential for large-scale data collection and aggregation.

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

### Participant interface

Participants can complete the Web-CDI on a variety of devices, including personal computers and tablets. Web-CDI can be administered on a smartphone, although the experience is not as ideal for the user due to the length of the survey. (As Web-CDI moves in the future to incorporate more short forms and adaptive forms, smartphone-responsive design will become a priority.) When a participant clicks a URL shared by a researcher, they are directed to a website displaying their own personal administration of the Web-CDI, regardless of whether the link was participant-specific or general-purpose. In some cases, they may be asked to read and accept a waiver of consent documentation, depending on whether the researcher has chosen to use that feature (see also Researcher

Interface below).

*Demographics.* The participant is next asked to provide demographic information about their family and any health conditions that might impact their child’s vocabulary development. Researchers can customize the presentation of these demographic questions in three ways. First, they can elect to show all of the demographics items on the landing page or to present the majority of these questions at the end of the instrument. This choice is provided because some pilot work in the United Kingdom indicated that answering questions regarding personal health information early in administration may deter participants from completing the instrument. Second, certain demographic questions can be asked at both the beginning and the end of the form to serve as validity checks, such that participants’ answers to redundant questions can be compared in order to screen for hasty or illegitimate completions. Third, researchers can tailor the questions to the societal and cultural context of their participants (e.g., country-specific education level descriptors, income categories, ethnicity definitions, etc.).

*Instructions.* After completing the first demographics page, participants are provided with instructions that are appropriate for either the Words & Gestures or Words & Sentences version (see Figure 1). At the top of the page are general instructions that inform participants that they should expect the study to take at least 30 minutes and that they should try to complete it in a quiet setting (e.g., while their child is sleeping). In addition, there are more detailed instructions for completing the vocabulary checklist. Unlike the traditional paper versions, instructions on how to properly choose responses are provided both in written and pictorial form. The pictorial instructions (Figure 1) aim to further increase caregivers’ understanding of how to complete the checklist. For example, these instructions clarify that the child’s understanding of a word requires them to have some understanding of the object that the word refers to or some aspect of the word’s meaning. In addition, caregivers are reassured that “child-like” forms (e.g., “raff” for “giraffe”) or family- or dialect-specific forms (e.g., “nana” for “grandma” are acceptable).

**Instructions:**

- This form can be filled any time 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:** Feb 15, 2020 6:02 PM



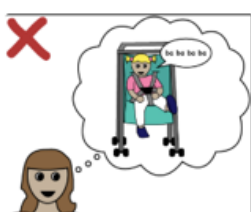
[Show waiver of documentation](#)  
Reach out to the Web-CDI Team!

[Save](#)

[Go back to Background Info](#)

[Table of Contents](#)

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.

Figure 1. Pictorial instructions in the Web-CDI Words and Sentences instrument.

Lastly, caregivers are reminded that the child should be able to produce the words “on their own” and that imitations are not acceptable. These general “rules of thumb” for completing the form should be familiar to researchers who are distributing the forms to parents so they can field any questions that may arise. While this is not possible for certain use-cases (e.g., collecting data via Facebook), these instructions should ideally also be reviewed either in writing (e.g., via email) or verbally (e.g., over the phone), so that these pictured instructions serve merely as a reminder to caregivers when completing the form.

*Completing the instrument.* The majority of the participant’s time in the study is spent completing the main sections of the instruments. As shown in Figure 2, on the American English Words and Gestures form, the vocabulary checklist portion of the form (396 items) asks parents to indicate whether their child can “understand” or “understand



**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 "spagetti") 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 "spagetti") 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"/> grrr        |
| <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.

and say” each word. Gesture communication and other early milestones are also assessed. In the American English Words and Sentences form, the vocabulary checklist (680 items) only asks parents to indicate which words their child “says”. Additional items assess children’s production of their three longest sentences, as well as morphological and syntactic development more broadly. All of these items are broken up across multiple screens for easier navigation through the form.

At the completion of the form, a graph is displayed illustrating the proportion of words from each semantic category that the child currently produces or understands. In addition, data from the norming studies are used to estimate the “hardest” (i.e., most advanced) word that the child currently understands or produces. This feedback to parents

is intended to provide parents with a fun “thank you” and is intentionally not designed to provide specific feedback about their child’s progress relative to other children or any normative standard. The closing page also reminds parents that their participation does not constitute a clinical evaluation and that they should contact their pediatrician or primary care physician if they have any concerns about their child’s development.

## Researcher interface

Table 1

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

| Study setting   | Default value     | Notes  |
|---|-------------------|--|
| Study name  | none              | NA   |
| Instrument  | none              | NA   |
| 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                 | NA   |
| 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. |

Table 1

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

| Study setting  | Default value                              | Notes  |
|--|--|--|
| 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.  |
| 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  | NA   |
| Would you like participants to be able to share their Web-CDI results via Facebook?                                | No   | NA   |

Table 1

*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 participants to answer the confirmation questions? | No                                  | Asks redundant demographic questions to serve as attention checks. |
| 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                                  | NA   |
| End message   | Standard<br>end-of-study<br>message | Can be changed to customize end-of-study message.                  |

One of the main goals of Web-CDI is to provide a unified CDI platform to the child language research community. To that end, researchers request an account by contacting a member of the CDI Advisory Board. Once they have registered an account they can create studies to distribute to participants. One rationale for this personalized registration process is that we ask that researchers allow fully anonymized data from their participants to be shared with the CDI Advisory Board, so that it can be added to Wordbank (<http://wordbank.stanford.edu/>; Frank, Braginsky, Yurovsky, & Marchman, 2017) and

shared with broader research community. However, there is an opt-out option if researchers do not wish to share their data, making data contribution voluntary.

A study in the context of the Web-CDI system is a set of individual administrations created by a researcher that share certain specifications. Table 1 gives an overview of the customizable features that are available at the study level in Web-CDI. These features are set when creating a study for the first time in Web-CDI using the “Create Study” tool, and most of the features can be updated continuously during data collection using the “Update Study” tool. While some of these features are only particularly relevant to specific use cases (e.g., longitudinal research and social media data collection, described below), others are relevant to all researchers using Web-CDI.

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

Researchers can download participant data in two formats. Both formatting options output a comma-separated values file with one row per participant; the full data option includes participant-by-item responses, and allows researchers to explore item-level trends, while the summary data option omits item-level data and only provides summary scores (e.g., total number of words understood/produced, percentile scores by age and gender).

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

*Individual recruitment.* One possible workflow using Web-CDI is to send unique study URLs to individual participants. Researchers do so by entering numerical participant

IDs or by auto-generating a specified quantity of participant IDs, each with its own unique study URL, using the “Add Participants” tool in the researcher dashboard. New participants can be added on a continual basis so that researchers can adjust the sample size of their study during data collection. Unique links generated for individual participants expire, by default, 14 days after creation, though the amount of days before link expiration is adjustable, which may be an important consideration for some researchers depending on their participant populations and specific project timelines. This workflow is most suitable for studies which pair the CDI with other measures, or when researchers contact specific participants from an existing database.

*Longitudinal studies.* Web-CDI also facilitates longitudinal study designs in which each participant completes multiple administrations. Researchers wishing to design longitudinal studies can do so by entering a list of meaningful participant IDs using the “Add Participants” tool in the researcher dashboard. If a certain participant ID is added multiple times, Web-CDI will create multiple unique study URLs in the study dashboard that have the same specified ID. In addition, when creating studies, researchers can select whether they would like the demographics information, vocabulary checklist, or no sections at all to be prefilled when a participant fills out a repeat administration of the instrument. Unless researchers are interested in cumulative vocabulary counts, it is strongly recommended that they do not use the option to pre-fill the vocabulary checklist portion of the instrument in longitudinal administrations as parents should complete the instrument at each time point independently.

*Social media and survey vendors.* Web-CDI contains several features designed to facilitate data collection from social media recruitment or through third-party crowd-sourcing applications and vendors (e.g., Amazon Mechanical Turk, Prolific). First, rather than creating unique survey links for each participant, researchers can also use a single, anonymous link. When a participant clicks the anonymous link, a new administration with a unique subject ID is created in the study dashboard. Additionally,

Web-CDI studies have several customizable features that are geared towards anonymous online data collection. For example, researchers can adjust the minimum amount of time a participant must take to fill out the survey before they are able to submit; with a longer minimum time to completion, researchers can encourage a more thorough completion of the survey. Researchers can also ask participants to verify that their information is accurate by checking a box at the end of the survey, and can opt to include certain demographic questions at both the beginning and end of the survey, using response consistency on these redundant items as a check of data quality.

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

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

## System Design

Web-CDI is constructed using open-source software. All of the vocabulary data collected in Web-CDI are stored in a standard MySQL relational database, managed using

Django and Python and hosted either by Amazon Web Services or by an European Union compliant server (see below). Individual researchers can download data from their studies through the researcher interface, and Web-CDI admins have access to the entire aggregate set of data from all studies run with Web-CDI. Website code is available in a GitHub repository <https://github.com/langcog/web-cdi>, where interested users can browse, make contributions, and request technical fixes.

## Data Privacy and GDPR Compliance

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

In the European Union (EU), research data collection and storage is governed by the Generalized Data Protection Regulation (GDPR) and its local instantiation in the legal

---

<sup>1</sup> 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 re-identification in some cases and if potential users have questions, we encourage them to consult with an Institutional Review Board.



system of the member states. Some of the questions on the demographic form contain information that may be considered sensitive (e.g., information about children’s developmental disorders), and in some cases, the possibility of linking this sensitive information to participant IDs exists, particularly when researchers draw on local databases that contain full names and addresses for recruitment and contacting. As a result, issues regarding GDPR compliance arise when transferring data outside the EU, namely to Amazon Web Services servers housed in the United States. Following GDPR regulations, these issues would make a data sharing agreement between data collectors and Amazon Web Services necessary. In addition, all administrators who can access the collected data would have to enter such an agreement, which needs updating whenever personnel changes occur. To overcome these hurdles and in consultation with data protection officers, we opted to exploit the local technical expertise and infrastructure to set up a sister site housed on GDPR-compliant servers, currently available at <http://webcdi.mpi.nl>. This site is updated synchronously with the main Web-CDI website to ensure a consistent user experience and access to the latest features and improvements. This site has been used in 135 successful administrations so far and is the main data collection tool for an ongoing norming study in the Netherlands. We are further actively advertising the option to use the European site to other labs who are following GDPR guidelines and are planning adaptations to multiple European languages, where copyright allows.

### Current Web-CDI Usage

One of the key benefits of Web-CDI use is that the system in effect becomes a centralized repository for standardized administrations of the CDI, contributing anonymized data (again, on an opt-out basis) to future research and norming efforts. In this section, we provide some preliminary analyses of the American English Web-CDI, demonstrating the potential of the Web-CDI system to provide a distributed platform for gathering large CDI datasets.

Table 2

*Exclusions from full WebCDI sample*

| <b>Exclusion</b>                | <b>WG</b>         | <b>% of full</b> | <b>WS</b>         | <b>% of full</b> |
|---------------------------------|-------------------|------------------|-------------------|------------------|
|                                 | <b>exclusions</b> | <b>WG sample</b> | <b>exclusions</b> | <b>WS sample</b> |
|                                 |                   | <b>excluded</b>  |                   | <b>excluded</b>  |
| 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    | 363               | 12.66%           | 236               | 6.57%            |
| System error in word tabulation | 1                 | 0.03%            | 4                 | 0.11%            |
| Total exclusions                | 1292              | 45%              | 1646              | 46%              |

At time of writing, researchers from 15 universities in the United States have collected over 5,000 administrations of the American English CDI using Web-CDI since it was launched in late 2017, with 2,868 administrations of the WG form before exclusions and 2,868 administrations of the WS form before exclusions. We excluded participants from the subsequent analyses based on a set of stringent criteria intended for the creation of future normative datasets. We excluded participants if it was not their first administration of the survey; if they were born prematurely or had a birthweight under 5.5 lbs ( $< 2.5$  kg); reported more than 16 hours of exposure to a language other than English per week on average (amounting to  $>10\%$  exposure to English); had serious vision impairments, hearing deficits or other developmental disorders or medical issues<sup>2</sup>; completed the survey

<sup>2</sup> 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.

unrealistically quickly (defined here as in fewer than 8.5 minutes)<sup>3</sup>; or were outside of the correct age range for the survey. The exclusion criteria we used were similar to those used in Fenson et al. (2007), who adopted stringent criteria to establish vocabulary norms that reflect typically developing children’s vocabulary trajectories. A complete breakdown of the number of participants excluded on each criterion is in Table 2. Of the completed WG forms, 1,292 were excluded, leading to a final WG sample size of 1,576 administrations, and 920 WS administrations were excluded, leading to a final WS sample size of 1,948.

### Demographic distribution and exclusions

Figure 3 shows the distribution of participant ethnicities as compared with previously reported numbers in a large scale norming study of the paper-based CDI form by Fenson et al. (2007). White participants still comprised nearly three quarters of the Web-CDI sample, while a higher proportion of participants report mixed ethnic identification as compared to the 2007 norms. Few participants identified as Hispanic/Latino: 6.5% of WG participants and 5.1% of WS participants reported Hispanic or Latino heritage. The low percentage of Hispanic/Latino participants was due in part to our exclusion of children with substantial exposure to languages other than English. Participants’ educational attainment level was similarly skewed. Over 80% of children in the Web-CDI sample came from families with college-educated mothers compared to 43% from the same group in the 2007 norms (Figure 3). Furthermore, less than 1 percent of participants in our families report a maternal education level less than a high school degree, compared to 7% from the same group in the 2007 norms. The overrepresentation of white Americans with high levels

---

<sup>3</sup> This timing criterion was chosen by authors B.D. and V.M. during recent online data collection as a lenient cutoff, i.e., one that errs on the side of including, rather than excluding, participants; on paper-based forms, caregivers are told the test generally takes 20-40 minutes. We noted that in early rounds of recent data collection, most participants who completed the survey in less than 8.5 minutes reported floor-level vocabulary scores regardless of age.

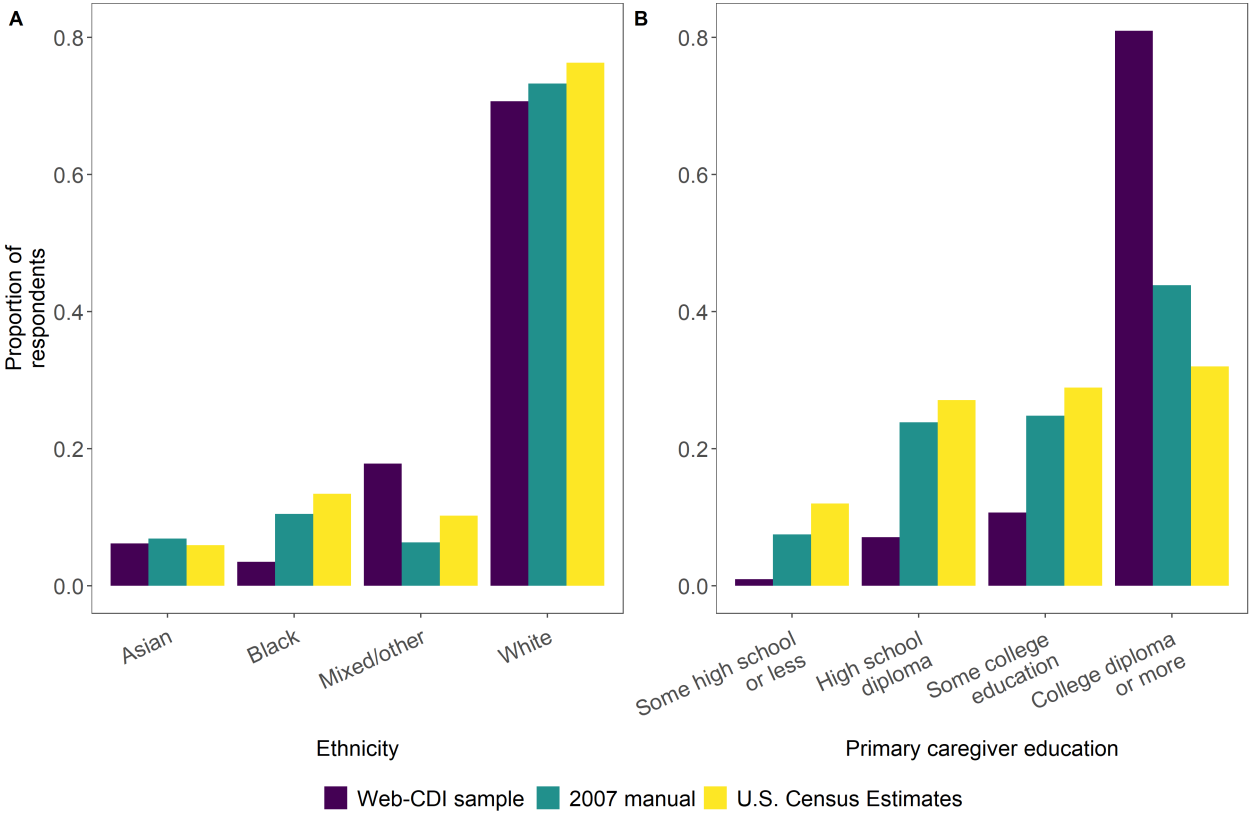
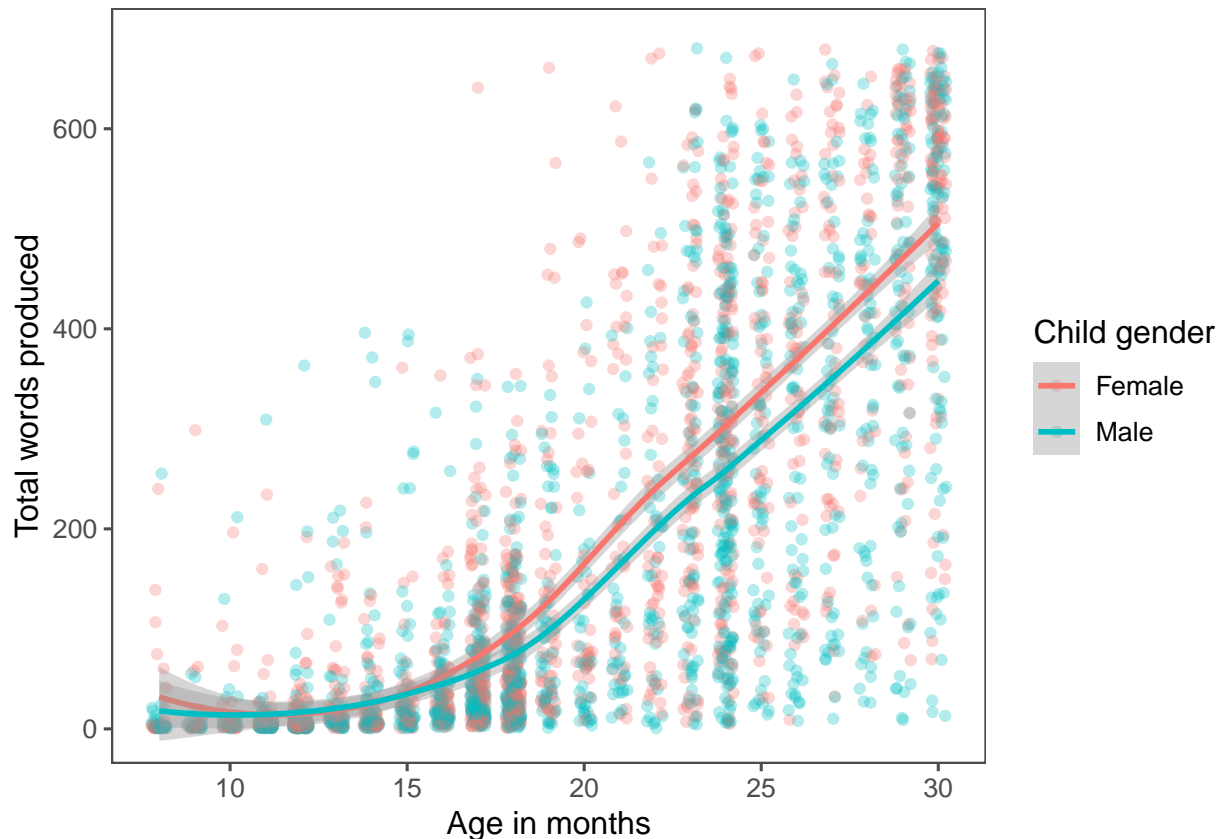


Figure 3. Proportion of respondents plotted by child race (A) and educational level of primary caregiver (B) from full Web-CDI sample to date (N = 3,524), compared with norming sample demographics from Fenson (2007). Latinx participants can be of any race.

of education attainment in this sample points to a general challenge encountered in vocabulary development research, which we return to when we detail our efforts to recruit more diverse participants.

### Results

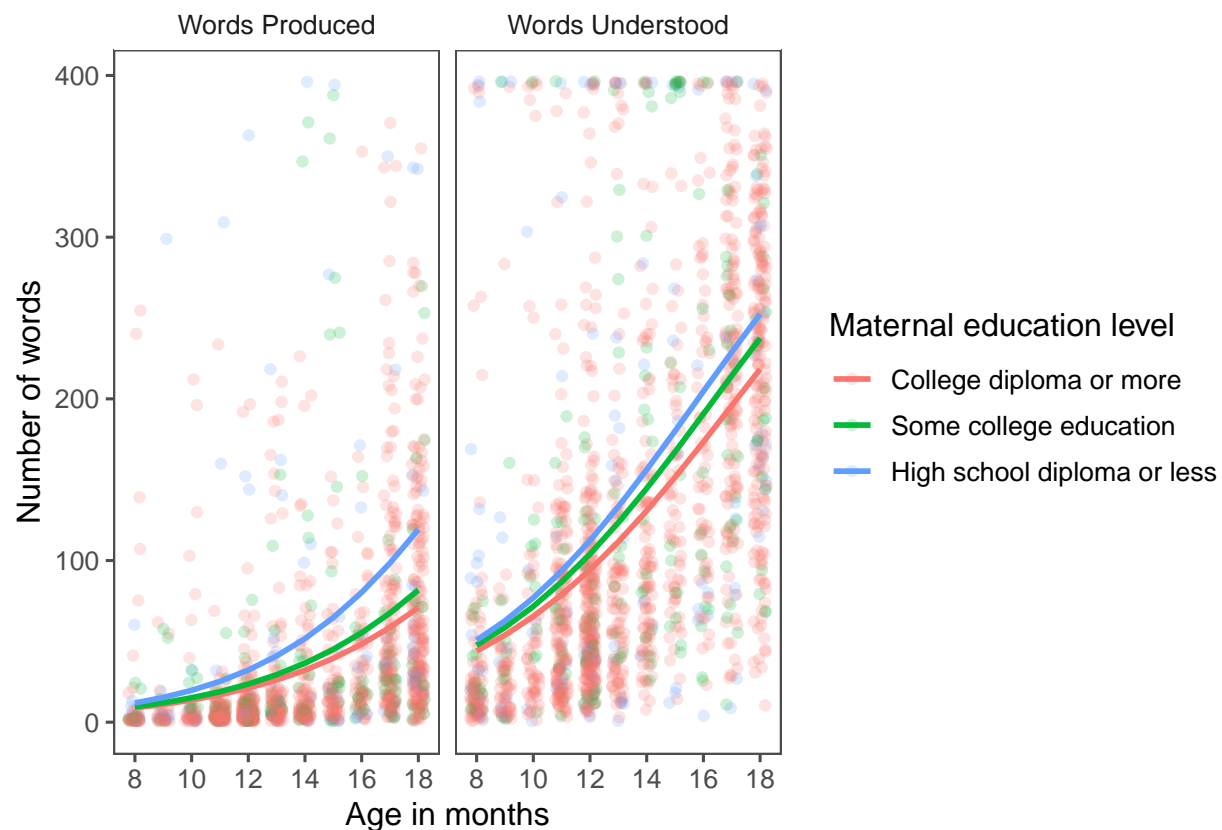
Although the CDI instruments include survey items intended to measure constructs other than vocabulary size, such as gesture, sentence production and grammar, we focus exclusively on the vocabulary measures here. Across both the WG and WS measures, our current Web-CDI sample shows greater reported vocabulary comprehension and production



*Figure 4.* Individual children’s vocabulary production scores from the entire Web-CDI sample plotted by children’s age and gender (both WG and WS,  $N = 3,513$ , with 1,674 girls). Line is a locally weighted regression with associated 95% confidence interval. Children with a different or no reported gender ( $N = 11$ ) are omitted here.

for older children. Moreover, data from both measures replicate a subtle but reliable pattern such that female children tend to have slightly larger vocabulary scores than male children across the period of childhood assessed in the CDI forms (Frank et al., 2021), though in these data this difference does not appear until around 18 months (Figure 4).

On the WG form, respondents’ reports of children’s vocabulary comprehension and production both increased with children’s age (Figure 5). We replicate overall patterns found by Feldman et al. (2000) in that, on both the “Words Understood” and “Words Produced” measures, vocabulary scores were slightly negatively correlated with primary



*Figure 5.* Individual children’s word production (left panel) and comprehension (right panel) scores 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”) as reported in the sample of Words and Gestures Web-CDI administrations collected as of November 2020 ( $N = 1,576$ ). Curves show generalized linear models fits.

caregivers’ education level, such that those parents without any college education reported higher vocabulary scores on both scales. A linear regression model with robust standard errors predicting comprehension scores with children’s age and primary caregivers’ education level (binned into categories of “High school diploma or less”, “Some college education” and “College diploma or more<sup>4</sup>”) as predictors shows main effects of both age

<sup>4</sup> “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.

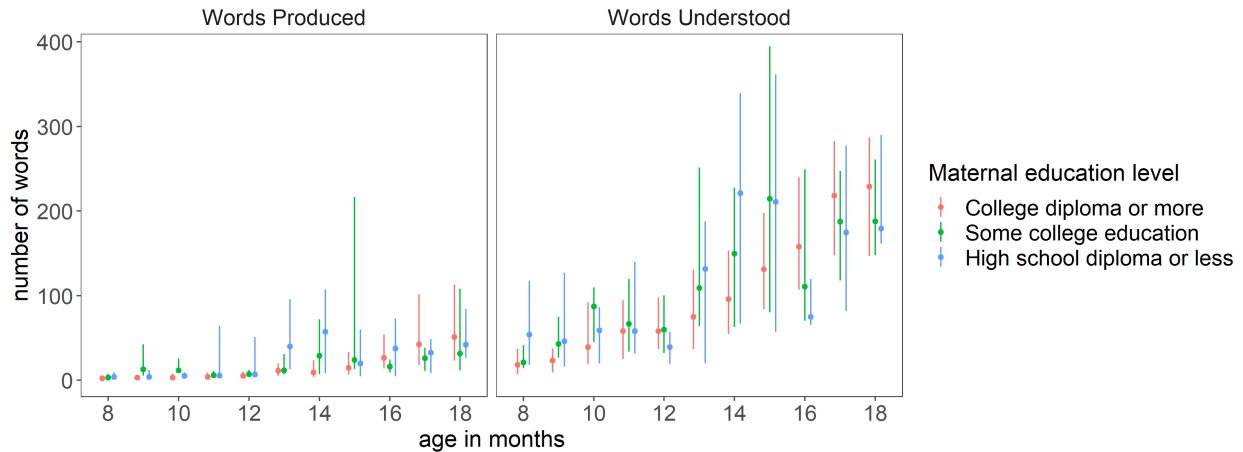
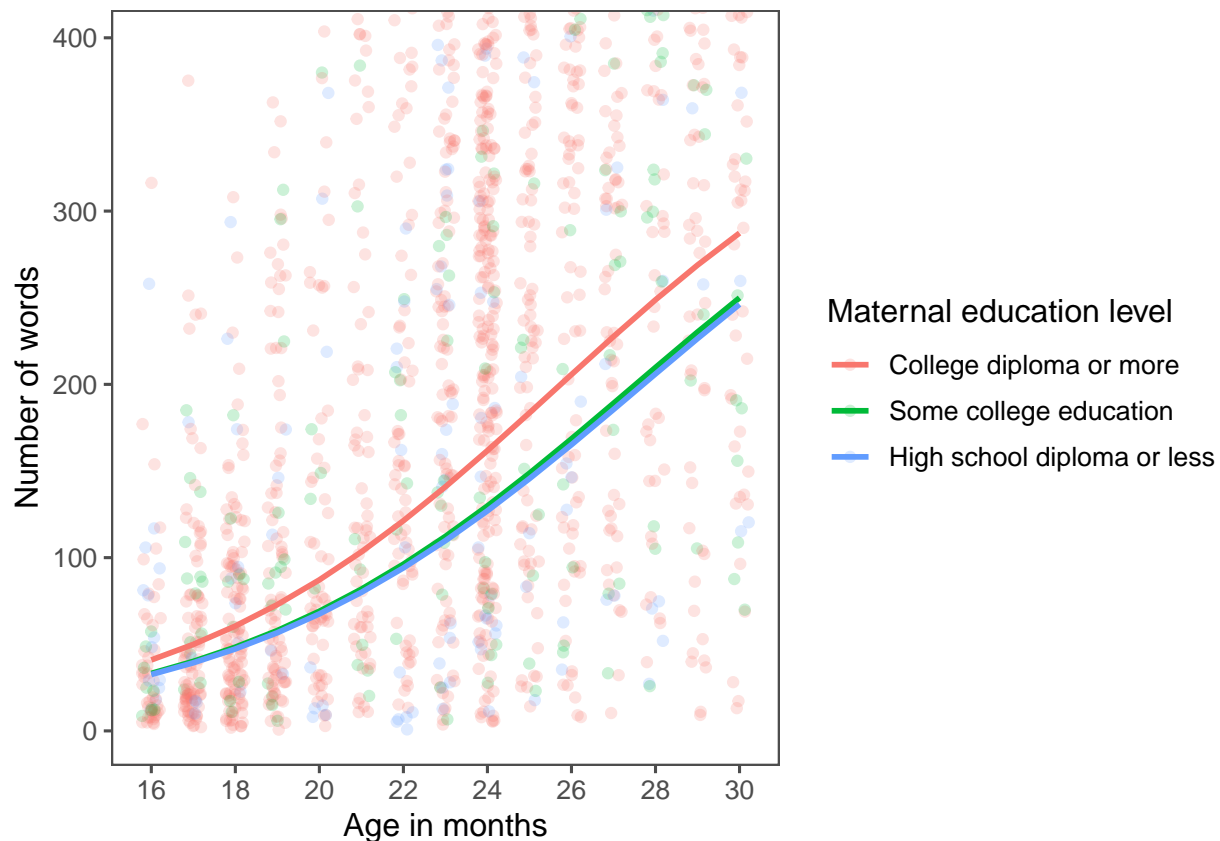


Figure 6. Median vocabulary production (left) and comprehension (right) scores by age on the WG form. Lines indicate span between first and third quartiles for each age.

( $\beta = 19.89$ ,  $p < 0.001$ ) and caregiver primary education ( $\beta_{highschool} = 29.59$ ,  $p = 0.01$ ).

Similarly, a linear regression model with robust standard errors predicting production scores by children's age and primary caregivers' education level shows main effects of age ( $\beta = 7.82$ ,  $p < 0.001$ ) and caregiver primary education ( $\beta_{highschool} = 28.86$ ,  $p = 0.002$ ).

The pattern of results seen in the WG sample is consistent with prior findings indicating that respondents with lower levels of education attainment report higher vocabulary comprehension and production on the CDI-WG form (Feldman et al., 2000; Fenson et al., 1994). Although caregivers with lower levels of education attainment report higher mean levels of vocabulary production and comprehension, median vocabulary scores (which are more robust to outliers) show no clear pattern of difference across primary



*Figure 7.* Individual children’s vocabulary production scores plotted by children’s age and maternal education level of primary caregiver education as reported in the sample of Words and Sentences Web-CDI administrations collected as of November 2020 ( $N = 1,948$ ). Curves show generalized linear models fits.

caregiver education levels (Figure 6). This discrepancy between the regression effects and a group-median analysis suggests that the regression effects described previously are driven in part by differential interpretation of the survey items, such that a few lower-SES caregivers are more liberal in reporting their children’s productive and comprehensive vocabularies, especially for the youngest children, driving up the mean scores for this demographic group.

Vocabulary production scores on the WS form show the expected pattern of increase with children’s age in months; in addition, scores replicate the trend reported in Feldman et al. (2000) and Frank et al. (2021) such that maternal education is positively associated



with children’s reported vocabulary size (Figure 7). Because representation of caregivers without a high school diploma is scarce ( $N = 18$  out of a sample of 1,948), interpretation of the data from this group is constrained. Nevertheless, as shown in Figure 7, a small but clear positive association between maternal education and vocabulary score exists such that college-educated caregivers report higher vocabulary scores than those of any other education level. The implications from these data converge with previous findings which indicate that parental education levels, often used as a metric of a family’s socioeconomic status, are related to children’s vocabulary size through early childhood.

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

### **Using Web-CDI to Collect Data from Diverse U.S.-based Communities**

Despite the large sample sizes we collected in the previous section, our current dataset from Web-CDI is, if anything, even more biased towards highly-educated White families than previous datasets. How can we recruit more diverse samples to remedy this issue? Here we discuss some potential routes forward. In this first effort we focus on collecting data from monolingual English-speaking families. While understanding that the performance of standard measurement tools like the CDI among multilinguals is of immense import to the field of vocabulary development research (Gonzalez et al., in prep; Floccia et al., 2018; De Houwer, 2019), we focused here only on monolingual development, because collecting data from multilingual populations introduces additional methodological considerations (e.g., how to measure exposures in each language) that are not the focus of our work here.

## Online data collection

Online recruitment methods, such as finding participants on platforms such as Amazon Mechanical Turk, Facebook and Prolific, represent one possible route towards assembling a large, diverse sample to take the Web-CDI. These methods allow researchers depart from their typical geographical recruitment area much more easily than with paper-and-pencil administration. However, these recruitment methods are, to our knowledge, largely untested with parent report measures of child language development. In a series of data collection efforts, we used Web-CDI as a tool to explore these different channels of recruitment.



Figure 8. Example Facebook advertisement in Phase 1 of recent data collection.

In our first phase of data collection, we ran advertisements on Facebook which were aimed at non-white families based on users' geographic locations (e.g., targeting cities

Table 3

*Exclusions from recent data collection using Facebook and Prolific.*

| <b>Exclusion</b>                | <b>WG</b>         | <b>% of full</b> | <b>WS</b>         | <b>% of full</b> |
|---------------------------------|-------------------|------------------|-------------------|------------------|
|                                 | <b>exclusions</b> | <b>WG sample</b> | <b>exclusions</b> | <b>WS sample</b> |
|                                 |                   | <b>excluded</b>  |                   | <b>excluded</b>  |
| 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    | 132               | 47.65%           | 122               | 40.40%           |
| System error in word tabulation | 0                 | 0.00%            | 0                 | 0.00%            |
| Total exclusions                | 162               | 58%              | 176               | 58%              |

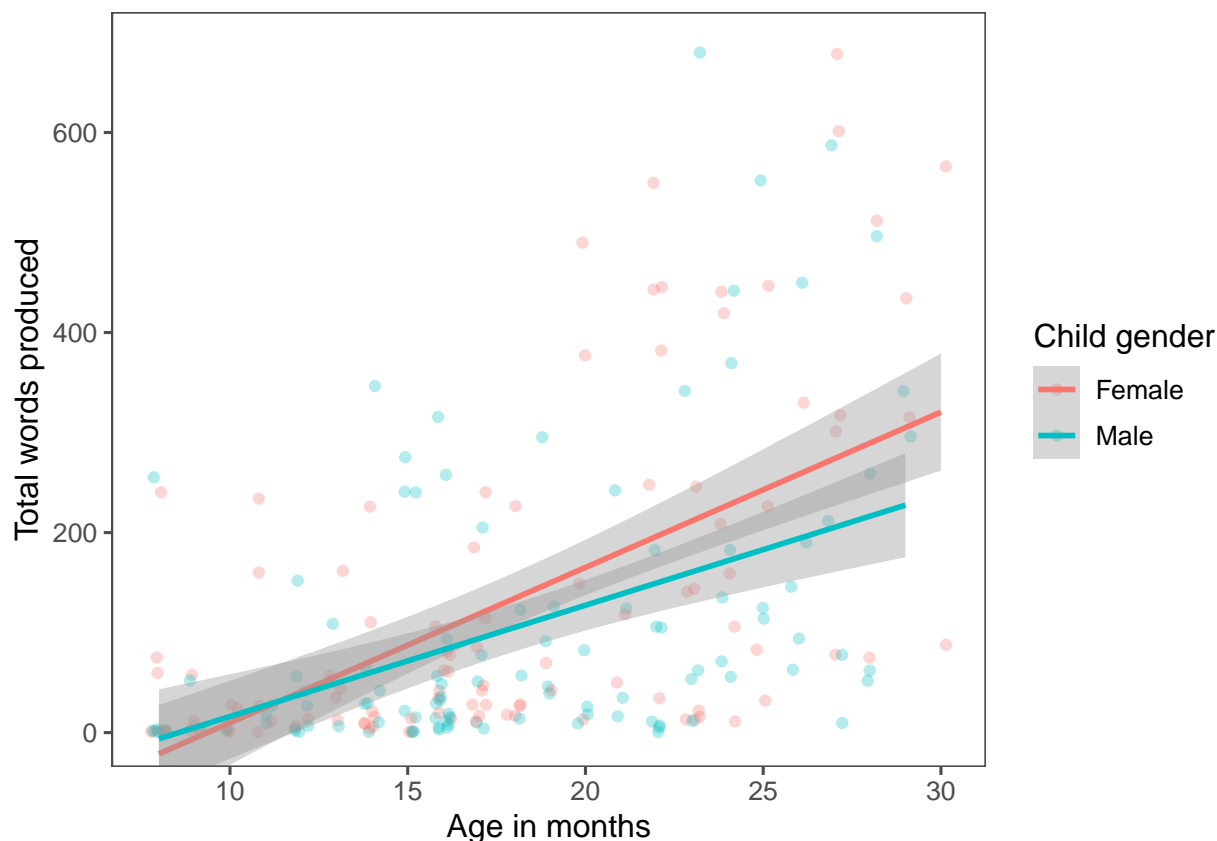
which have a higher than average representation of African Americans) or other profile features (e.g., ethnic identification, interest in parenthood-related topics). Advertisements consisted of an image of a child and a caption informing Facebook users of an opportunity to fill out a survey on their child’s language development and receive an Amazon gift card (Figure 8). Upon clicking the advertisement, participants were redirected to a unique administration of the Web-CDI, and they received \$5 upon completing the survey. This open-ended approach to recruitment offered several advantages, namely that a wide variety of potential participants from specific demographic backgrounds can be reached on Facebook. However, we also received many incomplete or otherwise unusable survey administrations, either from Facebook users who clicked the link and decide not to participate, or those who completed the survey in an extremely short period of time (over half of all completed administrations, Table 3).

In the second phase of our data collection efforts, we used the crowdsourcing survey vendor Prolific (<http://prolific.co>) in the hopes that some of the challenges encountered with Facebook recruitment would be addressed. Prolific allows researchers to create studies and post them to individuals who are in the platform’s participant database, each of whom is assigned a unique alphanumeric “Prolific ID.” Importantly, Prolific maintains detailed demographic information about participants, allowing researchers to specify whom they would like to complete their studies. Prolific further has a built-in compensation infrastructure that handles monetary payments to participants, eliminating the need to disburse gift cards through Web-CDI.

In the particular case of Web-CDI, the demographic information needed to determine whether an individual was eligible to complete our survey (e.g., has a child in the correct age range, lives in a monolingual household, etc.) was more specific than the information that Prolific collects about their participant base. We therefore used a brief pre-screening questionnaire to generate a list of participants who were eligible to participate, and subsequently advertised the Web-CDI survey to those participants. Given that we were interested only in reaching participants in the United States who were not white or who did not have a college diploma, our data collection efforts only yielded a sample that was small ( $N = 71$ ) but much more thoroughly screened than that which we could obtain on Facebook.

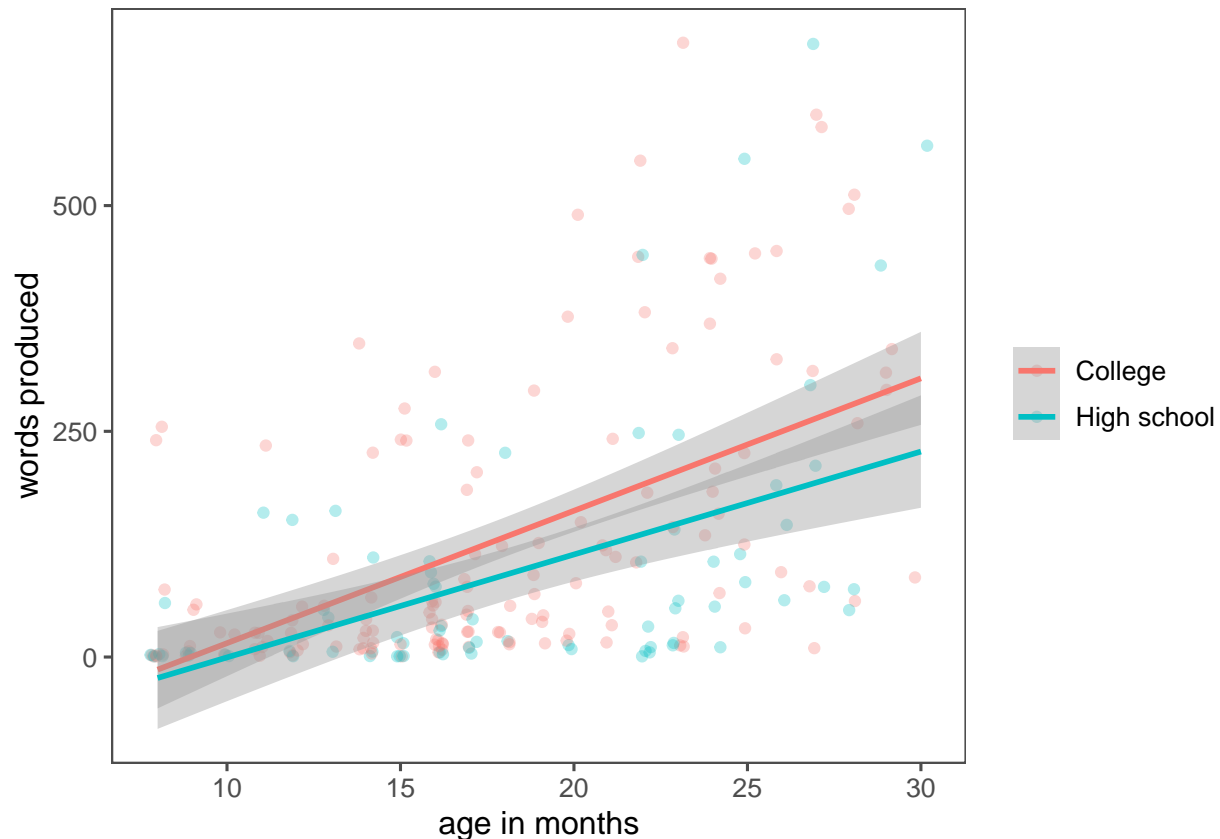
Across both phases (Facebook and Prolific recruitment), we used the same exclusion criteria as in the full Web-CDI sample to screen participants. A complete tally of all excluded participants is shown in Table 3. In both the WG and WS surveys, exclusion rates were high, amounting to 58% of participants who completed the survey. The high exclusions rates were notably driven by an accumulation of survey administrations which participants completed very quickly (in these analyses, defined as a completion taking less than 8.5 minutes). Many of the survey administrations excluded for fast completion had missing demographic information reported: Among WG participants excluded for too-fast

468 completions, 93% did not report ethnicity, and among WS participants excluded for the  
 469 same reason, 97% did not report ethnicity. After exclusions, full sample size was  $N = 115$   
 470 WG completions and  $N = 126$  completions.



*Figure 9.* Individual children’s vocabulary production scores from the entire Web-CDI sample plotted by children’s age and gender (both WG and WS,  $N = 238$ , with 116 girls). Lines are best linear fits with associated 95% confidence intervals. Children with a different or no reported gender ( $N = 3$ ) are omitted here.

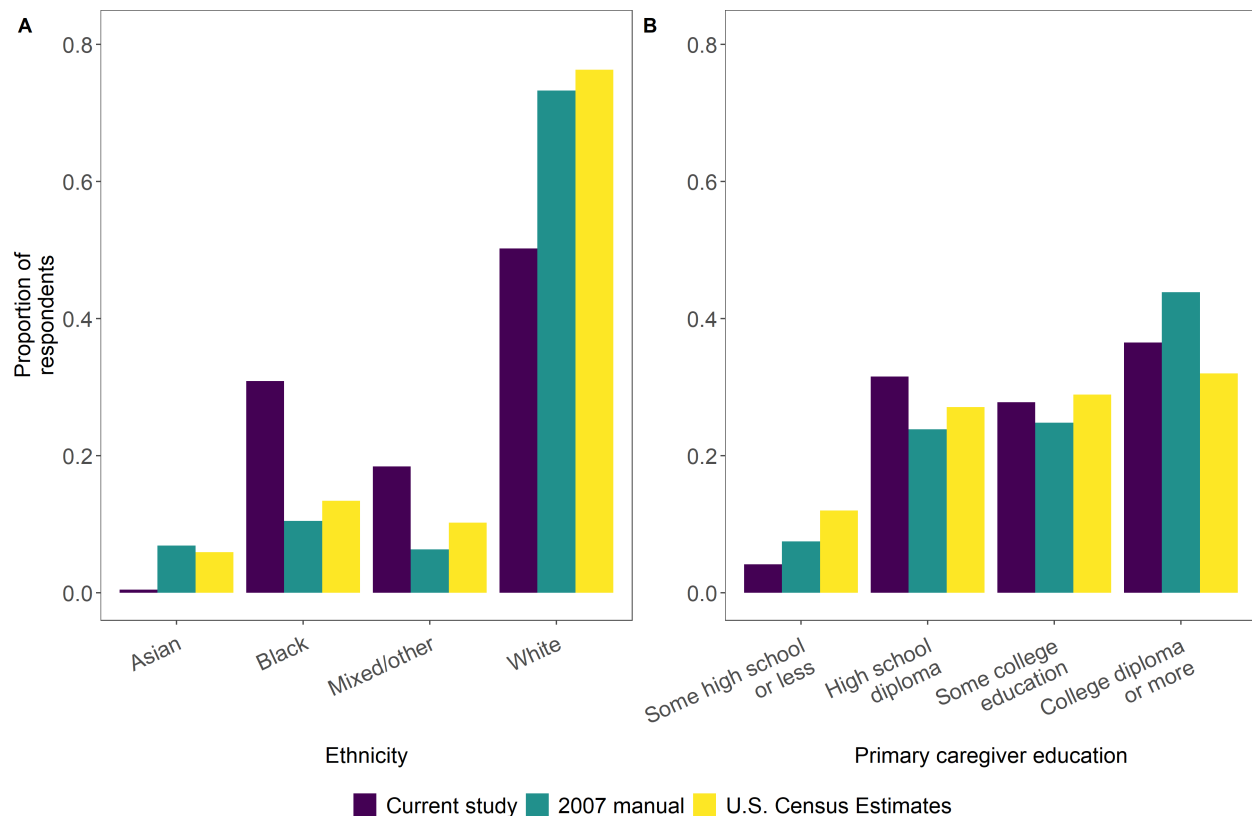
471 The results from our recent data collection efforts focused on lower-SES, non-white  
 472 participants show overall similar patterns to the full Web-CDI sample in several regards.  
 473 Word production scores from both the WG and WS administrations reflect growing  
 474 productive vocabulary across the second and third years, with a very small gender effect  
 475 such that female children’s vocabularies grow at a slightly faster rate than males’ (Figure



*Figure 10.* Individual children’s vocabulary production scores 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 ( $N = 241$ ). Lines show best linear fits and associated 95% confidence intervals.

9). The relationship between caregivers’ reported levels of maternal education and child’s vocabulary score is not as clear as it is in the full Web-CDI sample (Figure 10); however, children of college-educated parents show slightly faster vocabulary growth than do children of parents without any college degree. These patterns suggest that our data show similar general patterns to other CDI datasets with other populations (Frank et al., 2021).

Importantly, recent data collection efforts showed a substantial improvement in reaching non-white or less highly-educated participants. After exclusions, the Web-CDI data we collected through Facebook and Prolific have a higher proportion of non-white



*Figure 11.* Proportion of respondents plotted by child race (A) and educational level of primary caregiver (B) from recent data collection efforts aimed towards oversampling non-white, lower-SES families ( $N = 241$ ), compared with norming sample demographics from Fenson (2007). Latinx participants can be of any race.

participants than the overall Web-CDI sample and the norms established by Fenson et al. (2007) (Figure 11). Black participants in particular showed a marked increase in representation, from 10.5% in the 2007 norms to 30.9% in the recent sample, while the proportion of white participants decreased from 73.3% in the 2007 norms to 50.2% in the recent sample. Representation on the basis of families' reported primary caregiver education also improved (Figure 11). Participants with only a high school diploma accounted for 31.5% of the recent sample as compared to 23.8% in the 2007 norms, and representation of those with a college diploma or more education decreased from 43.8% in

the 2007 norms to 36.5% in the recent sample.

## Discussion

Taken together, these recent results indicate that Web-CDI could be a promising avenue through which to collect vocabulary development data in non-white, lower-SES communities when paired with online recruitment methods that yield legitimate, representative participant samples. These data do, however, convey clear limitations of our approach. Perhaps most conspicuously, more than half of completed administrations in this sample had to be excluded, in many cases because the information provided by participants appeared rushed or incomplete: over 40% of administrations were completed in fewer than 8.5 minutes, and of these quick completions, well over 90% were missing demographic information that is rarely missing in other administrations of the form. Determining the precise reasons for the high exclusion rate, and how (if at all) this (self-)selection may bias data reflecting demographic trends in vocabulary development, requires a more thorough assessment of who is submitting hastily-completed forms. This assessment is beyond the scope of the current study. However, all respondents who got to the end of the form were compensated regardless of how thoroughly they completed it, creating the possibility that some participants who clicked the anonymous link may not have been members of the population of interest, but rather were other individuals motivated by compensation.

Additionally, the exclusion rates described previously only provide information on those participants who did, at some point, submit a completed form, but many individuals clicked the advertisement link and did not subsequently continue on to complete the form. Without an in-depth exploration of who is clicking the link and why they might choose not to continue, we cannot draw conclusions about the representativeness of the current sample with regards to the communities we would like to include in our research. As such, a more thorough understanding of how users from different communities respond to various recruitment and sampling methods is needed in future work in order to draw conclusions



about demographic trends above and beyond those already established in the literature.

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

Finally, a significant limitation of the current data collection process is that many people in the population of interest - particularly lower-income families - do not have reliable internet access. Having participants complete the Web-CDI on a mobile device may alleviate some of the issues caused by differential access to Wi-Fi, since the vast majority of American adults own a smartphone (Pew Research Center, 2018 - NEED TO ADD MANUALLY TO REFS). Accordingly, improving Web-CDI's user experience on mobile platforms will be an important step towards ensuring that caregivers across the socioeconomic spectrum can easily complete the survey. For smartphone users on pay-as-you-go plans, who may be reluctant to use data to complete a study, a possible solution could be compensating participants for the amount of "internet time" they incurred completing the form.

## Conclusions

In this paper, we presented Web-CDI, a comprehensive online interface for researchers to measure children’s vocabulary by administering the MacArthur-Bates Communicative Development Inventory family of parent-report instruments. Web-CDI provides a convenient researcher management interface, built-in data privacy protections, and a variety of features designed to make both longitudinal and social-media sampling easy. To date, over 3,500 valid administrations of the WG and WS forms have been collected on Web-CDI from more than a dozen researchers in the United States after applying strict exclusion criteria derived from previous norming studies) (Fenson et al., 2007, 1994).

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

Beyond the goal of simply getting more data, we hope that Web-CDI can advance efforts to expand the reach of vocabulary research past convenience samples into diverse communities. A key question in the field of vocabulary development concerns the mechanisms through which sociodemographic variables, such as race, ethnicity, income and education are linked to group differences in vocabulary outcomes. Large, population-representative samples of vocabulary development data are needed to understand these mechanisms, but most research to date (including the full sample of

Web-CDI administrations) oversamples white participants and those with advanced levels of education.

We explored the use of Web-CDI as part of a potential strategy to collect data from non-white, lower-SES communities in two phases. Several overall patterns emerged from the resulting data which we expected: vocabulary scores grew with age, providing a basic validity check of the Web-CDI measure; females held a slight advantage in word learning over males; and children of parents with a college education showed slightly higher vocabulary scores. Nonetheless, the insights from these data, while aligned with past norming studies, are necessarily constrained by several features of our method. First, exclusion rates among data collected on Facebook were very high, well over 50%, mostly due to a large quantity of hasty completions. Second, a rigorous evaluation of the population-representativeness of those who were counted in the final sample was not feasible here.

Web-based data collection can capture useful information about vocabulary development from diverse communities, but future research will need to examine which sampling methods can yield accurate, population-representative data that can advance our understanding of the link between sociodemographic variation and variation in language outcomes.

## References

- Bates, E., & Goodman, J. C. (2001). On the inseparability of grammar and the lexicon: Evidence from acquisition.
- Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P., Reznick, J. S., . . . Hartung, J. (1994). Developmental and stylistic variation in the composition of early vocabulary. *J Child Lang*, 21(01), 85–123.
- Bornstein, M. H., & Putnick, D. L. (2012). Stability of language in childhood: A multiage, multidomain, multimeasure, and multisource study. *Developmental Psychology*, 48(2), 477.
- Caselli, N. K., Lieberman, A. M., & Pyers, J. E. (2020). The asl-cdi 2.0: An updated, normed adaptation of the macarthur bates communicative development inventory for american sign language. *Behavior Research Methods*, 1–14.
- 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).

Fenson, L., Pethick, S., Renda, C., Cox, J. L., Dale, P. S., & Reznick, J. S. (2000). Short-form versions of the macarthur communicative development inventories. *Applied Psycholinguistics*, 21(1), 95–116.

Floccia, C., Sambrook, T. D., Delle Luche, C., Kwok, R., Goslin, J., White, L., . . . others. (2018). Vocabulary of 2-Year-Olds learning english and an additional language: Norms and effects of linguistic distance. *Monographs of the Society for Research in Child Development*, 83(1), 1–135.

Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank: An open repository for developmental vocabulary data. *Journal of Child Language*, 44(3), 677–694. <https://doi.org/10.1017/S0305000916000209>

Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability and consistency in early language learning: The wordbank project*. MIT Press.

Kristoffersen, K. E., Simonsen, H. G., Bleses, D., Wehberg, S., Jørgensen, R. N., Eiesland, E. A., & Henriksen, L. Y. (2013). The use of the internet in collecting cdi data—an example from norway. *Journal of Child Language*, 40(03), 567–585.