

A thorough evaluation of the Language Environment Analysis (LENA) system

Alejandrina Cristia<sup>1</sup>, Marvin Lavechin<sup>1</sup>, Camila Scaff<sup>1</sup>, Melanie Soderstrom<sup>2</sup>, Caroline Rowland<sup>3</sup>, Okko Räsänen<sup>4,5</sup>, John Bunce<sup>2</sup>, & Erika Bergelson<sup>6</sup>

<sup>1</sup> Laboratoire de Sciences Cognitives et de Psycholinguistique, Département d'études cognitives, ENS, EHESS, CNRS, PSL University

<sup>2</sup> Department of Psychology, University of Manitoba, Canada

<sup>3</sup> Max Planck Institute for Psycholinguistics, Netherlands

<sup>4</sup> Unit of Computing Sciences, Tampere University, Finland

<sup>5</sup> Department of Signal Processing and Acoustics, Aalto University, Finland

<sup>6</sup> Psychology & Neuroscience, Duke University, Durham, North Carolina, USA

#### Author Note

Correspondence concerning this article should be addressed to Alejandrina Cristia, 29, rue d'Ulm, 75005 Paris, France. E-mail: alecristia@gmail.com

## Abstract

In the previous decade, dozens of studies involving thousands of children across several research disciplines have made use of a combined daylong audio-recorder and automated algorithmic analysis called the LENA<sup>®</sup> system, which aims to assess children’s language environment. While the system’s prevalence in the language acquisition domain is steadily growing, there are only scattered validation efforts, on only some of its key characteristics. Here, we assess the LENA<sup>®</sup> system’s accuracy across all of its key measures: speaker classification, adult word counts (AWC), child vocalization counts (CVC), and conversational turn counts (CTC). Our assessment is based on manual annotation of clips that have been randomly or periodically sampled out of daylong recordings, collected from (a) populations similar to the system’s original training data (North American English-learning children aged 3-36 months), (b) children learning another dialect of English (UK), and (c) slightly older children growing up in a different linguistic and socio-cultural setting (Tsimane’ learners in rural Bolivia). We find reasonably high accuracy in some measures (AWC, CVC), with more problematic levels of performance in others (CTC, precision and recall of male adults and other children). We find little difference in accuracy as a function of child age, dialect, or socio-cultural setting. Whether LENA<sup>®</sup> results are accurate enough for a given research, educational, or clinical application depends largely on the specifics at hand. We therefore conclude with a set of recommendations to help researchers make this determination for their goals.

*Keywords:* Speech technology; human transcription; English; Tsimane’; Reliability; Agreement; Method comparison; Measurement error; Child vocalization count; Adult word count; Conversational turn count; LENA

A thorough evaluation of the Language Environment Analysis (LENA) system

While nearly all humans eventually become competent users of their language(s), documenting the experiential context of early acquisition is crucial for both theoretical and applied reasons. Regarding theory, there are many open questions about what kinds of experiences and interactions are necessary, sufficient, or optimal for supporting language development. Moreover, the ability to accurately and quickly assess an infant’s state of development at a given point in time is of central importance for clinical purposes, both for children with known risks of language delays and disorders, and those who might not be identified based on risk factors. Reliable assessments are also crucial for measuring intervention efficacy.

One approach that has been making its way into the mainstream literature across basic and applied research on language and cognition relies on day-long recordings gathered with a LENA<sup>®</sup> audiorecorder (e.g., Greenwood, Thiemann-Bourque, Walker, Buzhardt, & Gilkerson, 2011; Gilkerson et al., 2017; Oller et al., 2010; VanDam & De Palma, 2018), and further analyzed using automated, closed-source algorithms. As we summarize below, this approach has many advantages, which may explain its expanding popularity. While over a hundred papers over the past two decades have used the output automatically provided by LENA<sup>®</sup>, only a handful include validity estimates (e.g., Weisleder & Fernald, 2013; d’Apice, Latham, & Stumm, 2019; Zimmerman et al., 2009), even fewer where validity estimation was the primary focus of the paper (e.g., Busch, Sangen, Vanpoucke, & Wieringen, 2018; Bulgarelli & Bergelson, 2019; Canault, Le Normand, Foudil, Loundon, & Thai-Van, 2016; Ganek & Eriks-Brophy, 2018; Lehet, Arjmandi, Dilley, Roy, & Houston, 2018). As a result, few studies report sufficient details about validation accuracy for one or more metrics, limiting the interpretability of the results of a meta-analytic assessment (cf. Cristia, Bulgarelli, & Bergelson, 2019). The work undertaken thus far also has some limitations, which are described further in the “Previous Validation” section below, and mentioned briefly next.

First, many validations or evaluations of LENA<sup>®</sup> omit analysis of the less-directly “relevant” categories of input like noise, silence, or overlap. Second, many LENA<sup>®</sup> evaluations rely on the output from LENA<sup>®</sup> as a starting point to either select portions of the file for manual annotation, or use its segmentation into talker turns or vocalizations as the unit of analysis rather than segmenting the audio from scratch. Both decisions can lead to inflated accuracy estimates. Here we endeavor to conduct an evaluation that is fully independent of the LENA<sup>®</sup> algorithms’ automated assessment, permitting a systematic, extensive, and independent evaluation of its key metrics.

In a nutshell, this paper reports on the performance of LENA<sup>®</sup> algorithms when compared to human annotations carried out in a set of clips extracted from daylong audiorecordings gathered from (a) a sample of children similar to the LENA<sup>®</sup> training set (i.e. infants and toddlers, growing up in North American English-speaking homes, and aged 3-36 months), (b) a group of similarly aged children learning a different dialect (UK English); and (c) slightly older children learning a different language in a very different sociocultural setting (Tsimane’-learning children in rural Bolivia).

**Brief introduction to LENA<sup>®</sup> products.** The LENA<sup>®</sup> system consists of hardware and software. The hardware component is a lightweight, sturdy, and easy-to-use recording device worn by a child in specialized clothing. The software is a suite of proprietary computer programs designed to provide automated quantitative analyses of the children’s auditory environment and their own vocalizations. The latter was developed over an extensive corpus of full day audio recordings gathered using their patented recording hardware (Xu, Yapanel, & Gray, 2009). The original dataset included over 65,000 hours of recording across over 300 American English-speaking families chosen for diversity in child age (1-42 months) and socio-economic status (Gilkerson & Richards, 2008). Half-hour selections from 309 recordings were transcribed and annotated for the purpose of developing the algorithm, and an additional 60 minutes from 70 additional recordings were transcribed

and annotated for testing the result (Gilkerson et al., 2008).

The resulting LENA<sup>®</sup> software takes as input a new audio recording and processes it incrementally in short windows, extracting a variety of acoustic features which are used to classify the audio stream into segments of at least 600 ms in length (or longer for some of the categories) using a Minimum Duration Gaussian Mixture Model (MDGMM; Xu et al., 2009). Silence may be included to “pad” segments to this minimum duration. The segments are classified according to a set of broad speaker and non-speaker classes. The speaker classes are: Male Adult, Female Adult, “Key” Child (i.e. the one wearing the recorder) and Other Child. The non-speaker classes are: Noise, Television (including any electronics), Overlap (speech overlapped with other speech or nonspeech sounds), and Silence (SIL). With the exception of Silence, these classifications are then passed through a further likelihood test between the original classification for a given segment and the Silence class, the result of which determines whether they are “Near” (high probability of being that class) or “Far” (low probability; i.e. they may be silence instead). Given the large number of acronyms and labels of various kinds, we provide a listing of relevant LENA<sup>®</sup> abbreviations in Table 1.

After this broad speaker classification step, Female or Male Adult “Near” segments (FAN and MAN) are further processed using an adaptation of the Sphinx Phone Decoder (Lamere et al., 2003) in order to form an automated estimate of the number of words in each segment (Adult Word Count, or AWC). Key Child (CHN) segments are further processed to sub-classify regions in them into vegetative noises, crying, and speech-like vocalizations. LENA<sup>®</sup> provides counts (child vocalization count, or CVC) and durations for this last speech-like sub-segment category. A further metric, Conversational Turn Counts (CTC), reflects the number of alternations between an adult and the key child (or vice versa), bounded by a maximum 5s of non-speech.

**Previous validation work.** A recent systematic review (Cristia et al., 2019) found 23 papers containing 28 studies that reported on the accuracy of the LENA<sup>®</sup> system’s labels

Table 1

*A partial listing of common LENA abbreviations and their meanings.*

Abbreviations	Meanings
FAN, MAN, CHN, CXN	Basic “meaningful speech” (near and clear speech) categories used by LENA for further processing: Female Adult Near, Male Adult Near, Key Child Near and Other Child Near categories respectively.
NON, TVN, OLN, SIL	Basic non-speech categories: Noise Near, Television Near, Overlap Near, Silence.
FAF, MAF, etc.	“Far” (low probability) versions of each Near category.
Key child	Child wearing recorder
AWC	Adult Word Count (estimated within FAN and MAN vocalizations)
CVC	Child Vocalization Count (estimated for non-cry, non-vegetative portions of CHN)
CTC	Conversational Turn Count (estimated for turns between FAN or MAN and CHN)

and/or derived metrics (AWC, CVC, CTC). They conclude that there are:

“reasonably good results [overall]: over 61% for recall and precision based on 11-12 non-independent studies; correlations for AWC mean  $r=.79$ , on  $n=11$ , with a mean RER [i.e., Relative Error Rate]=10% on  $n=11$ ; CVC mean  $r=.76$ ,  $n=5$ , with a mean RER=1% on  $n=5$ . The exception to this general trend towards good performance was CTC, with a mean  $r=.31$ ,  $n=5$ , RER=-64% on  $n=2$ .”

The systematic review also identified several limitations of previous validation work.

First, for the majority of included studies, the validation component was not fully evaluated by peer review. Even if the study may have appeared in a peer-reviewed journal, the validation in itself was often a secondary goal to support a different research objective, and therefore often lacked methodological details or even full results. For instance, Seidl et al. (2018) report on validation of LENA<sup>®</sup> labels among children at familial risk for autism in a one-paragraph appendix to the paper, which only mentions confusions between female adult and child. This leaves unclear whether confusions between Key child and any other category (Other child, Male adult, Silence, etc.) were ignored or considered to be errors. While this approach may be reasonable for a given study’s research goals, it has the undesirable side effect of creating the impression that LENA<sup>®</sup> metrics are widely validated, while in fact validation methods may not have been reported or evaluated in detail. Second, previous studies typically did not take silence, noise, or overlap into account in the reported confusion matrices or other accuracy measures, particularly within segments. That is, if a LENA<sup>®</sup> segment labeled “key child” contained one second of silence and two seconds of speech by the key child, the full three second clip may be tagged as “correct” though it was only 67% correct, leading to an overestimation of the accuracy of the “key child” label.

Third, a majority of previous validation studies used the LENA<sup>®</sup> output itself to select the sections that would be annotated for validation (in Cristia et al., 2019, this held for 14/25 studies that specified the method of selection). For instance, clips may have been selected for manual annotation on the basis of high AWC and/or CTC according to the algorithm. This unfortunately leads to biased sampling: Since LENA<sup>®</sup> only counts words within FAN and MAN segments and conversational turns involving FAN/MAN alternations with CHN in close proximity, high AWC or CTC can only occur in sections of the recording that are “clean” enough for the algorithm to parse; otherwise, most of the section would have been classified as overlap (OLN), which does not count towards AWC or CTC. This would tend to bias these reports toward a higher level of accuracy than would be obtained across the full recording.

Fourth, previous validation work has typically focused on a single corpus, participant population, age range, and language. As a result, although considerable variation in performance has sometimes been reported (e.g., Gilkerson et al., 2016; Canault et al., 2016) it is difficult to assess whether a numerical difference in accuracy found is significant, and if so, whether this is due to a difference in the way the corpus was constituted and annotated, rather than on how LENA<sup>®</sup> fares with that population, age range, and language.

**The present work.** We sought to assess the validity of the output provided by LENA<sup>®</sup> through an approach that complements the preceding literature. Specifically, we report an evaluation of all speech labels, also considering non-speech labels (notably silence, overlap, and TV, with limitations in our approach to be discussed below); as well as the system’s key derived metrics: adult word count (AWC), conversational turn count (CTC), and child vocalization count (CVC). We aim to address several of the limitations found in the body of previous work.

First, to maximally avoid potential bias in our annotations, we used random or periodic sampling (detailed below) to choose which sections of daylong recordings to annotate, and did not give annotators access to the LENA<sup>®</sup> output. Second, the fact that annotators did not have access to the LENA<sup>®</sup> segmentation allowed an assessment of the accuracy of the segmentation itself as well as categorical labeling. Specifically, LENA<sup>®</sup> and human levels were evaluated every 10 ms. This allows us to capture a much finer-grained representation of the auditory environment (i.e., if LENA<sup>®</sup> classified a 2 s audio segment as FAN, but .8 s of this was actually non-speech noise or a different talker, in our analysis LENA<sup>®</sup> would be credited only for the proportion that was correct).

Third, to gain traction on generalizability, rather than focusing on a single sample that either mirrors or diverges from LENA<sup>®</sup>’s original population, we included five corpora. Three corpora sampled from the same population, language, dialect, and age group the LENA<sup>®</sup> software was developed with. A fourth corpus was chosen to allow an extension to a different



dialect of English. The fifth corpus constituted an extension to a totally different recording condition (a rural setting, with large families and many children present, in a typologically different language). The age range also varies a great deal, and it is slightly higher in this last corpus. By and large, one could expect accuracy to decline in the sample of children who spoke a different English dialect compared to the three samples that matched better the data the LENA<sup>®</sup> software was developed with; and one could predict an even greater reduction in accuracy for the group that is learning a completely different language and which further mismatches in age (see other work on age- and language-mismatching samples, Busch et al., 2018; Canault et al., 2016).

Finally, the present study relies on a collaborative effort across several labs. The annotation pipeline was identical for four of the corpora, and conceptually comparable to the fifth (as detailed below). This allows us to more readily answer questions regarding differences in reliability as a function of e.g. child age and language. This approach also let us better infer the likelihood with which our results will generalize to other corpora, provided the annotation scheme is conceptually comparable.

## Methods

This paper was written using RMarkdown (Baumer, Cetinkaya-Rundel, Bray, Loi, & Horton, 2014) in R (Team & others, 2013) running on Rstudio (RStudio Team, 2019). It can be downloaded and reproduced using the data also available from the Open Science Framework, <https://osf.io/zdg6s>. These online Supplementary Materials also include a document with the full output of all models discussed here as well as additional analyses.

**Corpora.** The data for the evaluation comes from five different corpora, annotated in the context of two research projects. The largest one is the ACLEW project (Bergelson et al., 2017; Soderstrom et al., n.d.); in this paper we focus on four different corpora of child

daylong recordings that have been pooled together, sampled, and annotated in a coordinated manner. These four corpora are: the Bergelson corpus (“BER”) from US English families from the upstate New York area (Bergelson, 2016), the LuCiD Language 0–5 corpus (“L05”) consisting of English-speaking families from Northwest England (Rowland, Bidgood, Durrant, Peter, & Pine, 2018), the McDivitt and Winnipeg corpora (“SOD”) of Canadian English families (McDivitt & Soderstrom, 2016), and the Warlaumont corpus (“WAR”) of US English from Merced, California (Warlaumont, Pretzer, Walle, Mendoza, & Lopez, 2016). Some recordings in BER, and all recordings in SOD and WAR are available from HomeBank repository (VanDam et al., 2016). The second project contains a single corpus collected from Tsimane’ speaking families in Bolivia (“TSI”; Scaff, Stieglitz, Casillas, & Cristia, n.d.). Socioeconomic status varies both within and across corpora. Key properties of the five corpora are summarized in Table 2.

Table 2

*Key properties of the five corpora*

Corpus	Children	Clips	Clip duration (seconds)	Mean Age [range] (months)	Location
WAR	10	150	120	6.3 [3-9]	Western US
BER	10	150	120	11.2 [7-17]	Northeast US
SOD	9	150	120	12.3 [2-32]	Western Canada
L05	10	150	120	20 [11-31]	Northwest England
TSI	13	272	60	34 [15-58]	Northern Bolivia

Despite these differences, all five corpora consists of long (4–16 hour) recordings collected as children wear a LENA<sup>®</sup> recorder in a LENA<sup>®</sup> vest throughout a normal day and/or night. For the four ACLEW corpora, out of the 106 recorded participants, daylong recordings from 10 infants were originally chosen from each corpus were for manual annotation, selected to represent a diversity of ages (0–36 months) and socio-economic

contexts. In the SOD corpus, sensitive information was found in one of the files, and thus one child needed to be excluded. The tenth day for this corpus was a second day by one of the 9 included children. From each daylong file, fifteen 2-minute non-overlapping sections of audio (with a 5-minute context window) were randomly sampled from the entire daylong timeline for manual annotation. In total, this led to 20 hours of audio, and 4.6 hours of annotated speech/vocalizations (collapsing across all speaker categories).

The TSI corpus consisted of 1 or 2 recordings from 13 children, out of the 25 children recorded from field work that year; the other 12 had been recorded using other devices (not the LENA<sup>®</sup> hardware). From these files, 1-minute segments were sampled in a periodic fashion. That is, for each recording, we skipped the first 33 minutes to allow the family to acclimate to the recorder, and then extracted 1 minute of audio (with a 5-minute context window) every 60 minutes, until the end of the recording was reached. This resulted in a total of 4.5 hours of audio, and 0.7 hours of speech (collapsing across all speaker categories).

We chose to sample 1 or 2 minutes at a time (TSI, and ACLEW corpora, respectively) because conversations are likely to be bursty (Goh & Barabási, 2008). That is, it is likely the case that speech is not produced at a periodic rate (e.g., one phrase every 20 seconds), but rather it occurs in bursts (a conversation is followed by a long period of silence between the conversational partners, followed by another bout of conversation, perhaps with different interlocutors, followed by silence, and so on). In this context, imagine that you sample a 5-second stretch. If you find speech in that stretch, then it is likely you have by chance fallen on a conversation bout; if you do not find speech, then you have likely found a silence bout. If you were to extend that selection out to several minutes, then it is likely that you will simply add more material from the same type (i.e. conversation bout or silence bout). As a result, any sampling method that favors medium-sized stretches (5-15 minutes) will tend to end up with samples that are internally homogeneous (throughout the 5-15 minutes, there is a conversation, or there is silence throughout). If smaller clips are sampled out, this

heterogeneity is still captured, but (keeping the total length of audio extracted fixed) the number of clips that can be extracted is larger, thus likely increasing the likelihood that results will generalize to a new section of the audio.

In the 5 corpora, the 1- or 2-min samples were annotated for all hearable utterance boundaries and talker ID. In ACLEW corpora talker IDs reflected unique individual talkers, but were coded in such a way to readily allow mapping onto LENA<sup>®</sup>s talker categories, e.g. key child, other child 1, female adult 1, female adult 2 (cf. Casillas et al., 2017; Bergelson et al., 2019 for the general annotation protocol; Soderstrom et al., n.d., for an introduction to the databases). The ACLEW datasets also had other coding levels which will not be discussed here. In the TSI corpus, only the key child and one female adult whose voice recurred throughout the day were individually identified, with all other talkers being classified on the basis of broad age and sex into male adult, female adult, and other children.

**Processing.** Several different time units are needed to clarify how each metric is calculated (see Figure 1). Clips refer to the 1- or 2-minute samples extracted from recordings (TSI corpus and ACLEW corpora, respectively). This is the basic unit at which child vocalization counts and conversational turn counts can be established. In addition, since most previous work evaluating adult word counts did so at the clip level, we do so here as well.

The other metrics require a more detailed explanation, conveyed graphically in Figure 1. The stretch of time that has been assigned to a speech or non-speech class by LENA<sup>®</sup> is a *segment*. In one clip, there may be just one long segment (e.g., the whole clip has been assigned to Silence by LENA<sup>®</sup>); or there may be more (e.g., the first 5 seconds are attributed to the key child, then there is a 50-second Silence segment, and the final 5 seconds are attributed to a Female Adult). In the LENA<sup>®</sup> system’s automated analysis, only one of these categories may be active at a given point in time. In contrast, colloquially, “utterance” or “vocalization” refers to stretches of speech detected by humans and assigned to different

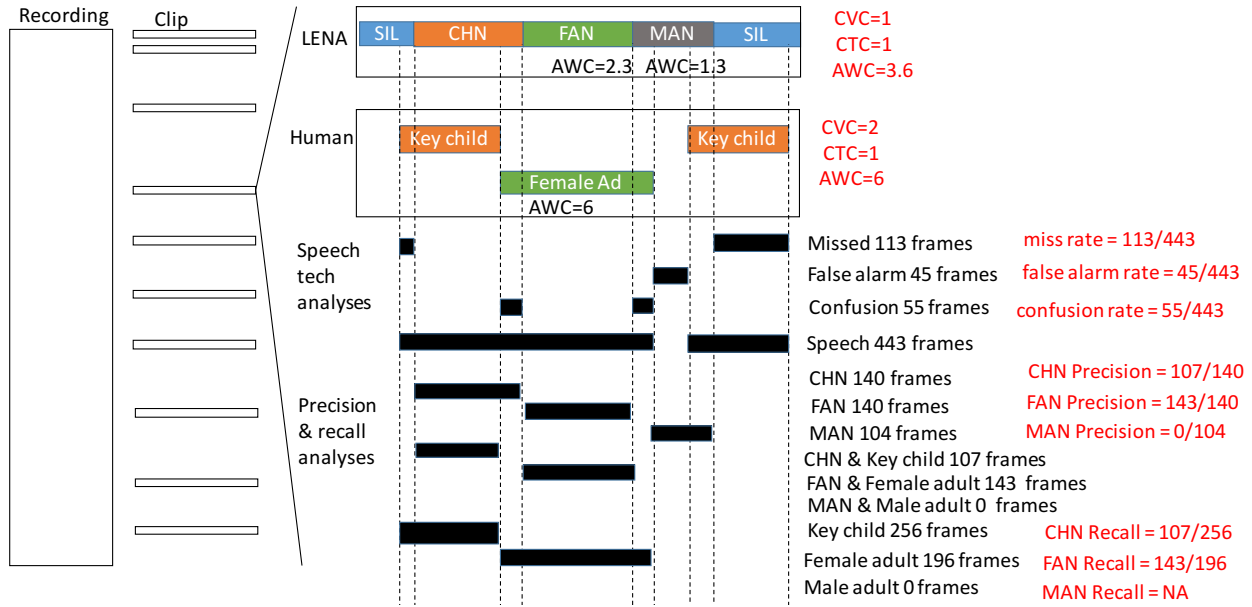


Figure 1. Levels at which performance is evaluated. Notice that there are multiple clips extracted from each recording; each clip can have zero or more segments; frames (10ms) are not shown because they would be too small in this scale. Adult Word Count (AWC), Child Vocalization Count (CVC), and Child Turntaking Count (CTC) are calculated at the level of the 1- or 2-minute long audio extracts (clips). Misses, false alarms, confusions as well as class precision and recall depend on 10-ms frames, and are totalled both at the level of individual clips and over the full audio extracts.

talkers. To be clear: in what follows, clips may have zero or more utterances. Unlike in the LENA<sup>®</sup> system, however, in our analysis a given point in time may be associated with multiple speakers.

Given that there need not be a one-to-one correspondence between LENA<sup>®</sup> segments and human utterances, we need to define smaller time units that can be used to check for classification agreement. In this paper, we use 10 ms *frames*. This is the basic time unit used for all classification accuracy estimations, which are introduced in more detail in the next subsection.

**LENA<sup>®</sup> classification accuracy.** Our first goal was to establish LENA<sup>®</sup> talker tag accuracy, particularly for the four broad LENA<sup>®</sup> talker categories (key child, other child, female adult, male adult; or CHN, CXN, FAN, MAN), while taking into account other categories (with some limitations on their interpretation clarified below). We calculated accuracy in two complementary ways. First, we used three frame-based standard metrics of speech and talker segmentation to allow direct comparison with other systems in the speech technology literature (False Alarm Rate, Miss Rate, Confusion Rate). We also use Identification Error Rate, which is derived by summing the first three metrics; together these provide a stringent and standard test of accuracy. Second, we used frame-based precision and recall of each category to provide an intuitive representation of the error patterns shown by this system.

***Speech and talker segmentation metrics.*** The original coding was converted using custom-written python scripts into a standard adaptation of the “Rich Transcription Time Mark” (rttm) format (Ryant et al., 2019), which indicates, for each vocalization or segment, its start time, duration, and speaker. This representation was used in pyannote.metrics (Bredin, 2017) to compute four standard identification metrics: rate of false alarm for speech, rate of misses for speech, rate of confusion between talkers, and the derived identification error rate (IDER). These are calculated with the following formulas at the level of each clip, where FA (false alarm) is the number of frames during which there is no talk according to the human annotator but during which LENA<sup>®</sup> found some talk; M (miss) is the number of frames during which there is talk according to the human annotator but during which LENA<sup>®</sup> found no talk; C (confusion) is the number of frames correctly classified by LENA<sup>®</sup> as containing talk, but whose voice type has not been correctly identified (when the LENA<sup>®</sup> model recognizes female adult speech where there is male adult speech for instance), and T is the total number of frames that contain talk according to the human annotation:

- **False Alarm rate** =  $FA/T$  (T=Total # of frames that contain talk),

- Miss rate =  $M/T$ ,
- Confusion rate =  $C/T$ ,
- **ID**entification **E**rror **R**ate (IDER) =  $(FA+M+C)/T$

In the human annotation, there is no class representing overlapping speech as such. For the sake of completeness and full comparison with the LENA<sup>®</sup> model, if two or more different speech sources were active at the same time according to the human annotators, these frames have been mapped to the class “overlap” post hoc. This allows us to compare this Overlap class to the LENA<sup>®</sup> system’s OLN (and, for the precision/recall analysis introduced next, OLF). Therefore, in the most complete analysis, the confusion rate is computed based on the human-LENA<sup>®</sup> matches in Table 3.

Table 3

*Correspondances between LENA and our human annotation tags for each talker type.*

*Additional analyses remove one or both of the last two rows. \*Electronic voices were only annotated in the ACLEW dataset. N.B. Although some Tsimane’ families listen to the radio, radio speech was not annotated in the TSI corpus.*

Talker	LENA	Human
Key Child	CHN	CHI
Other Child	CXN	OCH
Female Adult	FAN	FA
Male Adult	MAN	MA
Electronics	TVN*	E*
Overlap	OLN	OL

However, our overlap category is not defined identically to the LENA<sup>®</sup> overlap category. For LENA<sup>®</sup>, overlap between any two categories is labeled OLN – i.e., Noise + TV would be counted towards overlap as would FAN+FAN; whereas for us, only overlap between two

talker categories (e.g., key child and female adult, noise was not coded) counts as overlap. Similarly, the TVN LENA<sup>®</sup> class is not equivalent to the electronic speech tag in the ACLEW coding, because the former also includes music, singing, crowd noise and any other sound coming from a TV or another electronic source, whereas the latter only includes speech from an electronic source. Therefore, additional analyses map these classes onto “Other” post hoc, so as to not penalize confusions involving them.

***Precision and recall.*** This evaluation looks in more detail at the pattern of errors, by assessing how LENA<sup>®</sup> and human annotators agreed and disagreed. In both precision and recall, the numerator is the intersection between a LENA<sup>®</sup> tag and a human tag (e.g., the number of frames that LENA<sup>®</sup> classified as CHN and the annotator classified as Key child). The denominator differs: To calculate precision, we divide that number by the total number of frames attributed to a category by LENA<sup>®</sup>, whereas for recall, we divide by the total number of frames attributed to a category by the human annotator.

***Agreement.*** When two or more annotators provide data on the same classification, one can calculate agreement. We report on Cohen’s  $\kappa$  as a measure of the extent to which LENA<sup>®</sup> and human annotators coincide in their labeling.

***CVC and CTC evaluation.*** From the human annotation, each vocalization by the key child counted towards the total Child Vocalization Count (CVC) for a given clip if and only if the vocalization had been annotated as being linguistic (canonical or non-canonical in the ACLEW notation). For the Conversational Turn Count (CTC), a sequence of key child and any adult (or vice versa) within 5 seconds counted towards the clip total CTC. The Pearson correlation across LENA<sup>®</sup> and human estimations was then calculated.

Users may also wish to interpret the number of vocalizations or turns found by LENA<sup>®</sup>. Therefore, it is important to also bear in mind errors, error rates, and absolute error rates. Despite the similarity in their names, these three metrics provide different information. We



define *error* as follows: given a LENA<sup>®</sup> estimate, how close is the human-generated value. This is calculated as  $NL - NH$ , where  $NL$  is the number according to LENA<sup>®</sup> and  $NH$  is the number according to humans; this is done separately for each clip. By averaging across clips, we then get an idea of the bias towards overestimation (if this number is positive) or underestimation (if this difference is negative).

In contrast to *error*, *error rate* computes this bias in relation to the actual number of vocalizations tagged by the human coder:  $(NL - NH) / NL$ . For instance, imagine that we find that LENA<sup>®</sup> errs by 10 vocalizations according to the average error; this means that, on average across short clips like the ones used here, the numbers by LENA<sup>®</sup> would be off by 10 vocalizations. By using the error rate, we can check whether this seemingly small difference is indeed small relative to the actual number found. That is, an error of 10 vocalizations would be less problematic if there were 100 vocalizations on average (in which case LENA<sup>®</sup> would be just 10% off) than if there were 10 (LENA<sup>®</sup> would be doubling the number of vocalizations). As with error, the sign of this difference indicates whether LENA<sup>®</sup> tends to over- or under-estimate these counts.

Finally, the *absolute error rate* is calculated with the formula  $\text{abs}(NL - NH) / NL$ , where  $\text{abs}$  indicates absolute value. As a result, it cannot be used to assess systematic under- or over-estimation biases, but rather gives an idea of how accurate the estimates are at the clip level (statistically speaking). To convey this intuitively, one could find an error of 0 together with an error rate of 0 because half of the samples are -100 vocalizations off (for the error) or -100% off (for the error rates), with the other half behaving in the exact opposite fashion. The *absolute error rate* then avoids this kind of cancellation by taking the polarity (+/-) of the error out of play.

**AWC evaluation.** For the AWC portion of this evaluation, we could only use transcriptions from the four ACLEW corpora, since the TSI corpus has not been transcribed (and thus lacks word counts). Annotators for the four ACLEW corpora were proficient in the

language spoken in the daylong recording, and transcribed all adult speech based using canonical lexical forms (e.g. “wanna”, not “want to”) in keeping with minCHAT format (MacWhinney, 2017).

Reference adult word counts were determined by counting all unambiguously transcribed words spoken by adult talkers. This was achieved by first discarding all non-lexical transcript entries such as non-linguistic communicative sounds, paralinguistic markers, and markers indicating incomprehensible speech. In addition, all utterances from the key child and other children were omitted from the Adult Word Count(AWC). The remaining orthographic entries separated by whitespaces were then counted as gold standard target words for LENA<sup>®</sup> to detect.

The 1- or 2-minute clips sampled for manual annotation were not guaranteed to perfectly align with LENA<sup>®</sup> segments (i.e. talker onsets and offsets), posing a potential issue for comparing LENA<sup>®</sup> Adult Word Counts relative to the human annotated word count. Of all LENA<sup>®</sup> segments found within the extracted clips, 14% straddled a clip boundary (i.e., the segment began before the clip started; or it ended after the extracted clip ended). To match LENA<sup>®</sup> AWCs with the annotated word counts, words from these straddling LENA<sup>®</sup> segments were included proportionally. That is, if 10% of the duration of a LENA<sup>®</sup> segment fell within a clip, 10% of the LENA<sup>®</sup> AWC estimate for that segment was included in the LENA<sup>®</sup> word count estimate for that clip. AWC was evaluated using Pearson correlations and error analyses, similarly to CVC and CTC.

## Results

Before starting, we provide some general observations based on the manual human annotations. The “Other” category (meaning no speech, potentially silence but also non-human noise) was extremely common, constituting 71% of the 10 ms frames. In fact,

Table 4

*Number of frames, percentage of frames,  
and number of minutes attributed to each  
category by the human annotators.*

	Frames	Percentage	Minutes
CHI	588,236	7	98
FEM	891,717	10	149
MAL	234,199	3	39
OCH	262,702	3	44
OVL	271,636	3	45
ELE	218,535	3	36
Other	6,112,975	71	1,019

30% of the 1-2 minute clips contained no speech by any of the speaker types (according to the human annotators). As for speakers, female adults made up 10% of the frames, the child contributed to 7%, and male adult voices, other child voices, and electronic voices were only found in 3% of the frames each. Overlap made up the remaining 3% of frames. The following consequences ensue. If frame-based accuracy is sought, a system that classifies every frame as Other would be 71% correct. This is of course not the kind of system we'd generally want, but it indicates that systems well adapted to this kind of recording should tend to have low false alarm rates, being very conservative as to when there is speech. If the system does say there is speech, then a safe guess is that this speech comes from female adults, who provide a great majority of the speech, nearly 1.5 times as much as the key child and 2 times more than other children or male adults. In fact, given that speech by male adults and other children is relatively rare, a system that makes a lot of mistakes in these categories may still have a good global performance, because males and other children jointly accounted for only 6% of the frames.

**LENA<sup>®</sup> classification accuracy: False alarms, misses, confusion.** Our first analysis is based on standard speech technology metrics, which establish errors relative to speech quantity. That is, if 10 frames are wrong in a file where there are 100 frames with speech, this is a much smaller problem than if 10 frames are wrong in a file where there is 1 frame with speech. In other words, these metrics should be considered *error rate* metrics.

However, one problem emerges when there is no speech whatsoever in a given file. In the speech technology literature, this is never discussed, because most researchers working on this are basing their analyses on files that have been selected to contain speech (e.g., recorded in a meeting, or during a phone conversation). We still wanted to take into account clips with no speech because this was central for our research goals: We need systems that can deal well with long stretches of Other (i.e., non-speech or silence), because we want to measure in an unbiased manner how much speech (and silence!) children hear.

Unfortunately, in the 30% of clips that had no speech whatsoever, the false alarm, miss, and confusion rates are all undefined, because the denominator is zero. To be able to take clips with no speech into account, we defined the following rules. First, if a clip had no speech according to the human annotator, while LENA said there was speech, then the false alarm rate was 100%, and the miss and confusion rates were zero. Second, if on the contrary, both the human annotator and LENA said there was no speech, then all the error rates were zero. Finally, in some cases there was just a little speech; in this case, the denominator was very small, and therefore the ratio for these two metrics ended up being a very large number. To be maximally informative, we report results in three ways: (1) *weighted by speech*: Overall false alarm, miss, and confusion rates over all clips together, thus giving more weight to clips with more speech; (2) *equal weight per clip*: means across clips, which represent central tendency when giving equal weight to clips with more versus less or no speech; and (3) *accounting for potential outliers*: since means are not robust to outliers, we also report the median across all clips.

As mentioned briefly above, there were, *a priori*, several ways of analyzing the data:

- one-on-one mapping of all the categories, as in Table 1;
- excluding TV as a speaker category, since it is conceptually not identical to the electronic voices detected by ACLEW human annotators; in this case, the TV labels would be mapped post hoc to “Other”, as would the electronic voices in the human annotation;
- omitting OLN as a speaker category, since it is not conceptually identical to the overlap derived from humans’ annotating different speaker categories; in this case, LENA<sup>®</sup>’s OLN labels would be mapped post hoc to “Other”, as would the regions of overlap in the human annotation.<sup>1</sup>

The analysis that yields the best LENA<sup>®</sup> performance (Table 5) focuses on the speaker categories while mapping electronic voices and overlap in the human annotation onto Other, so that the categories considered in the human annotation are FEM, MAL, CHI, OCH, alongside using only CHN, FAN, MAN, and CXN as speakers in the LENA<sup>®</sup> annotation, (with all “far” categories, TVN, and OLN all mapped onto Other; see Tables 1 and 3). Calculated in this way, LENA<sup>®</sup>’s *false alarm* rate (i.e., tagging a speech category when there was none) and *confusion* rate (i.e., providing the wrong label) were lowest. Notably, however, the *miss* rate (i.e., the system returns a judgment that no sound label is activated) was double that found with the other analysis alternatives. See Table 5.

In the second-best performing case, electronic voices in the human annotation are still mapped onto Other but Overlap is not, so that the human categories considered were CHI, FEM, MAL, OCH, and overlap; and the LENA<sup>®</sup> categories considered were CHN, FAN,

---

<sup>1</sup>Please note that the “near” and “far” versions of the categories cannot be used in the present analysis, which requires a direct one-to-one matching across the system labels and the gold labels. The division between “near” and “far” in LENA<sup>®</sup> is done by comparing the fit to a given category and silence with a likelihood ratio; such a computation was not possible for human annotators.

Table 5

*False Alarm Rate (FAR), Miss Rate (MR), Confusion Rate, and total Identification Error Rate (IDER, sum of the medians of the other three categories), as a function of which categories are considered. Speakers indicates that only speaker categories are considered (all others are mapped onto SIL); + Electronic that also electronic was scored; + Overlap that electronic and overlap in both human and LENA annotation were also scored.*

	Overall				Mean				Median			
	FAR	MR	CR	IDER	FAR	MR	CR	IDER	FAR	MR	CR	IDER
Speakers	13	56	11	79	26	39	NA	73	6	39	NA	73
+ Electronic	44	24	38	106	86	20	NA	132	20	12	NA	88
+ Overlap	58	22	42	122	126	17	NA	172	30	9	NA	98

MAN, CXN, and OLN (with all “far” classes and TVN mapped onto Other). Finally, performance was worst when included the electronic voices segmented by human annotators, and LENA<sup>®</sup>’s TVN speaker categories in the evaluation (rather than mapping them all to Other), such that the human categories considered were CHI, FEM, MAL, OCH, overlap, and electronic; and the LENA<sup>®</sup> categories considered were CHN, FAN, MAN, CXN, OLN, and TVN. It is likely that these differences are partially due to OLN and TVN not being defined similarly across the system and human annotators.

**LENA<sup>®</sup> classification accuracy: Precision and recall.** By now, we have established that the best performance emerges when “far” labels such as CHF and OLF are mapped onto Other, as are TVN/ELE and OLN/OVL. False alarm, miss, and confusion rates are informative but may be insufficient for our readers for two reasons. First, these metrics give more importance to correctly classifying segments as speech versus non-speech (false alarms + misses) than confusing talkers (confusion). Second, many LENA<sup>®</sup> users are

particularly interested in sections labeled speech, especially that of adults and the key child, rather than sections labeled as non-speech. The metrics reported thus far do not give more importance to adults and the key child, and they do not give us insight into the patterns of error made by the system.

We therefore turn to precision and recall. Looking at precision of speech categories is crucial for users who interpret the LENA<sup>®</sup> system’s estimated quantity of adult speech or key child speech, as low precision means that some of what LENA<sup>®</sup> called e.g. key child was not in fact the key child, and thus it is providing overestimates. Looking at recall may be most interesting for users who intend to employ LENA<sup>®</sup> as a first-pass annotation: the lower the recall, the more is missed by the system and thus cannot be retrieved (because the system labeled it as something else, which will not be inspected given the original filter).

This subsection shows confusion matrices, containing information on precision and recall, for each key category. For this analysis, we collapsed over all human annotations that contained overlap between two speakers into a category called “overlap”. Please remember that this category is not defined the same way as the LENA<sup>®</sup> overlap category. For LENA<sup>®</sup>, overlap between any two categories falls within overlap – i.e., CHN+TV would be counted towards overlap; whereas for the manual annotations, only overlap between two talker categories (e.g., key child and female adult) counts as overlap.

***LENA<sup>®</sup> classification accuracy: Precision.*** We start by explaining how to interpret one cell in Figure 2: Focus on the cross of the human category (i.e., row) FEM and the LENA<sup>®</sup> category (i.e., column) FAN; when LENA<sup>®</sup> tagged a given frame as FAN, this corresponded to a frame tagged as being a female adult by the human 60% of the time. This category, as mentioned above, was the most common speaker category in the audio, so that over 287k frames (representing 60% of the frames tagged as FAN by LENA<sup>®</sup>) were tagged as a female adult by both the human and LENA<sup>®</sup>. The remaining 40% of frames that LENA<sup>®</sup> tagged as FAN were actually other categories according to our human coders: 18% were

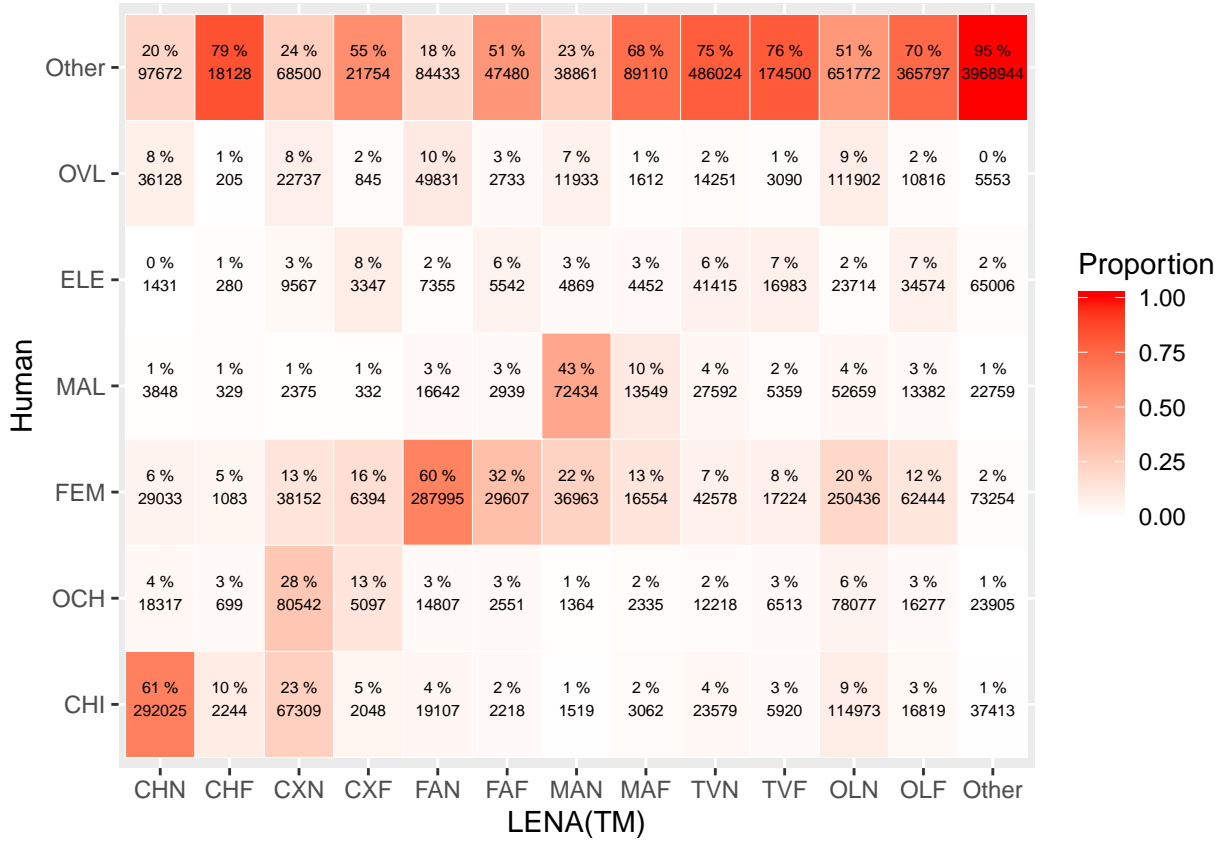


Figure 2. Precision: Confusion matrix between LENA (x axis) and human annotations (y axis). In each cell, the top number indicates the percentage of all frames in that LENA category (column) that are labeled as a given class by the human (row); cells in a given column add up to 100%. The number below indicates number of frames in that intersection of LENA and human classes.

Other (meaning no speakers were talking), 10% were in regions of overlap between speakers or between a speaker and an electronic voice, and 12% were confusions with other speaker tags. Inspection of the rest of the confusion matrix shows that, other than Other, FAN and CHN are the LENA<sup>®</sup> tags with the greatest precision.

Indeed, precision for CHN is almost identical, at 61%; thus, over half of the frames labeled as the key child are, in fact, the key child. The majority of the frames that LENA<sup>®</sup> incorrectly tagged as being the key child are actually Other (that is, silence or more generally



lack of speech) according to the human annotator (%), with the remaining errors being due to confusion with other categories. About 6% of them are actually a female adult; 4% are another child, and 8% are regions of overlap across speakers, according to our human coders.

Lower precisions are found for MAN (43%) and CXN(28%). The pattern of confusion is somewhat different from the other two categories we looked at, due to greater confusion with the other label within the same age class. That is, 22% of the frames LENA<sup>®</sup> tagged as being MAN actually corresponded to female adult speech according to the human annotation. It was also not uncommon to find a CXN tag for a frame human listeners identified as a female adult (13%), but even more confusions involved the key child (28%). In a nutshell, this suggests increased caution before undertaking any analyses that rely on the precision of MAN and CXN, since most of what is being tagged with these talker codes by LENA<sup>®</sup> is other speakers or Other (i.e. silence, absence of speech).

Another observation is that the “far” tags of the speaker categories do tend to more frequently correspond to what humans tagged as Other (absence of speech; 67%) than the “near” tags (35%), and thus it is reasonable to exclude them from consideration for most purposes.

The relatively high proportion of near LENA<sup>®</sup> tags that correspond to regions that humans labeled as Other could be partially due to the fact that the LENA<sup>®</sup> system, in order to process a daylong recording quickly, does not make judgments on short frames independently, but rather imposes a minimum duration for all speaker categories, padding with silence in order to achieve it. Thus, any key child utterance that is shorter than .6 s will contain as much silence as needed to achieve this minimum (and more for the other talker categories). Our system of annotation, whereby human annotators had no access whatsoever to the LENA<sup>®</sup> tags, puts us in an ideal situation to assess the impact of this design decision. That is, any manual annotation that starts from the LENA<sup>®</sup> segmentation would likely bias the human annotator to ignore such interstitial silences to a greater extent than if they have

no access to the LENA<sup>®</sup> tags. We inspected how often this padding by the LENA<sup>®</sup> system occurred and found that it was quite common: About half of the child linguistic and non-linguistic vocalizations tagged in any given clip were shorter than 600 milliseconds long, and thus would have been padded by LENA<sup>®</sup> with silence automatically.

These precision analyses shed light on the extent to which the LENA<sup>®</sup> tagged segments contain what the speaker tag name indicates, relative to human coders. We now move on to *recall*, which indicates a complementary perspective: how much of the original annotations attributed to a given class was captured by the corresponding LENA<sup>®</sup> class.

***LENA<sup>®</sup> classification accuracy: Recall.*** Again, we start with an example to facilitate the interpretation of Figure 3. As seen at the intersection of human CHI (last row) and LENA<sup>®</sup> CHN (first column), the best performance for a talker category for recall is CHN: 50% of the frames humans tagged as being uttered by the key child were captured by the LENA<sup>®</sup> under the CHN tag. Among the remainder of what humans labeled as the key child, 11% was captured by the LENA<sup>®</sup> system’s CXN category and 20% by its OLN tag, with the rest spread across several categories.

This result suggests that an analysis pipeline that uses the LENA<sup>®</sup> system to capture the key child’s vocalizations by extracting only CHN regions will get half of the key child’s speech. If additional manual human vetting is occurring in the pipeline, researchers may find it fruitful to include segments labeled as CXN, since this category actually contains a further 11% of the key child’s speech. Moreover, as we saw above, 28% of the CXN LENA<sup>®</sup> tags corresponds to the key child, which means that human coders re-coding CXN regions could filter out the 72% that do not, if finding key child speech were a top priority.

Many researchers also use the LENA<sup>®</sup> as a first pass to capture female adult speech through the FAN label. Only 32% of the female adult speech can be captured this way. Unlike the case of the key child, missed female speech is classified into many of the other

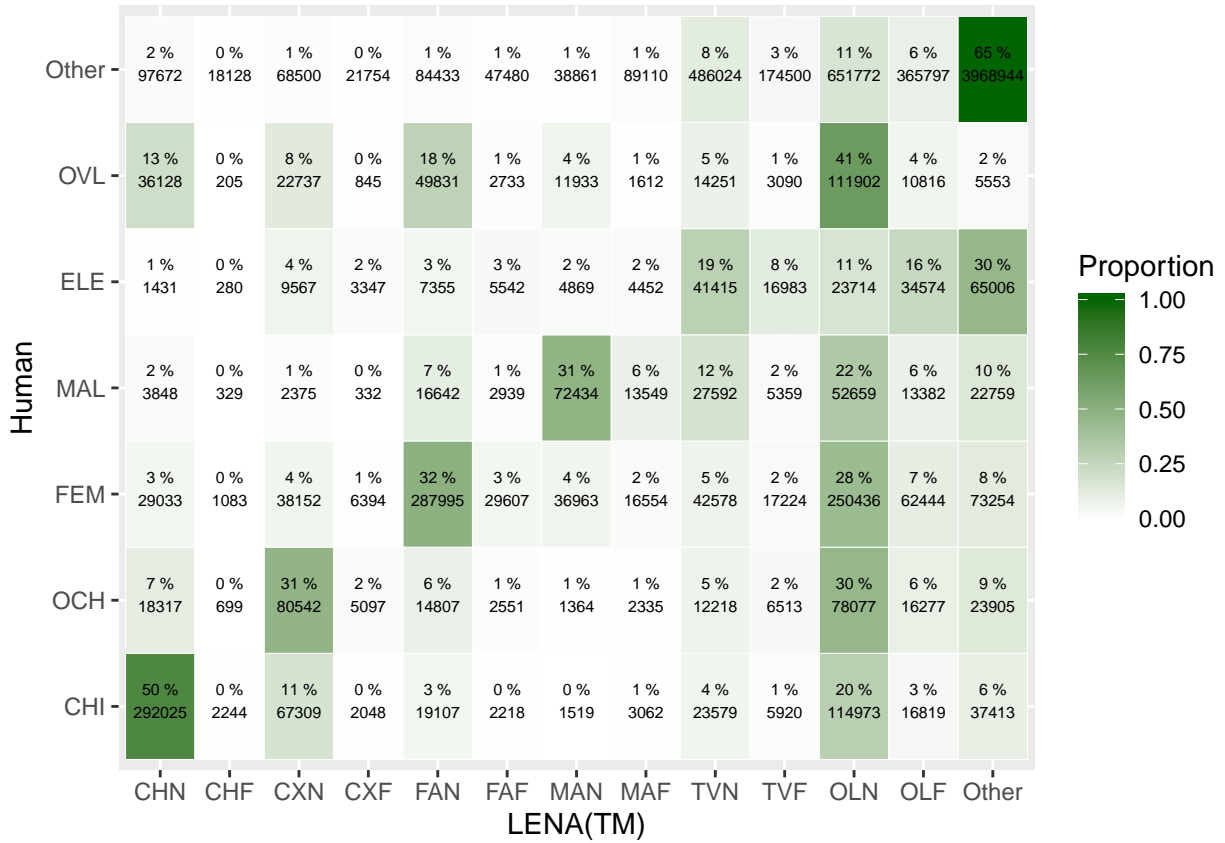


Figure 3. Recall: Confusion matrix between LENA (x axis) and human annotations (y axis). In each cell, the top number indicates the percentage of all frames that a human labeled as a given class (row) which were recovered in a given LENA category (column); cells in a given row add up to 100%. The number below indicates number of frames in that intersection of LENA and human classes.

categories, and thus there may not exist an easy solution (i.e., one would have to pull out all examples of many other categories to get at least half of the original female adult). However, if the goal is to capture as much of the female speech as possible, a reasonable solution would be to include OLN regions, since these capture a further 28% of the original female adult speech and, out of the OLN tags, 20% are indeed female adults (meaning that if human annotators are re-coding these regions to find further female adult speech, they'd need to filter out 80% of the segments, on average).

For the remaining two speakers (MAL, OCH), recall averaged 31%, meaning that a third of male adult and other child speech is being captured by LENA<sup>®</sup>. In fact, most of these speakers’ contributions are being tagged by LENA<sup>®</sup> as OLN (mean across MAN and CXN 22%) or TV (mean across MAN and CXN is 12%), although the remaining sizable proportion of misses is actually distributed across many categories.

Finally, as with precision, the “far” categories show worse performance than the “near” ones. It is worth noting that it is always the case that a higher percentage of frames is captured by the near rather than the far labels. For instance, out of all frames attributed to the key child by the human annotator, 50% were picked up by the LENA<sup>®</sup> CHN label whereas essentially 0% were picked up by the LENA<sup>®</sup> CHF label. This result provides further support that when sampling LENA<sup>®</sup> daylong files using the LENA<sup>®</sup> software, users likely need not take the “far” categories into account.

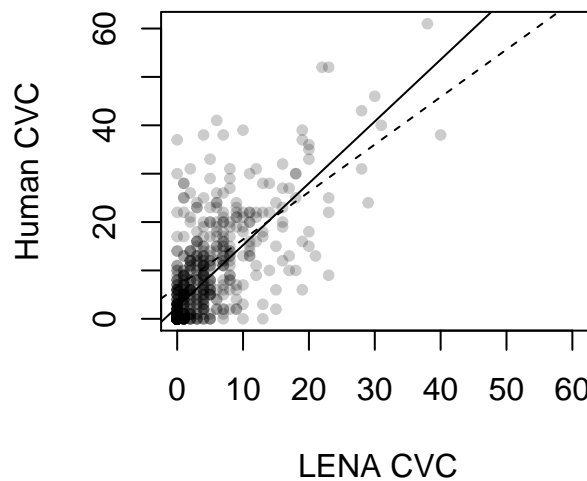
***LENA<sup>®</sup> classification accuracy: Agreement using Cohen’s  $\kappa$ .*** Given results above suggesting that our coding of electronics may not have coincided with the LENA<sup>®</sup> system’s, and that “far” categories are inaccurate, in this analysis we only consider the following labels for LENA<sup>®</sup>: (all others are collapsed into an Other category); and the following labels for human annotators: (all others are collapsed into an Other category). This analysis revealed a Cohen’s  $\kappa$  estimated at  $K(8580000) = 0.44$ , weighted  $\kappa$  estimated at  $K(8580000) = 0.46$ .

**Child Vocalization Counts (CVC) accuracy.** We extracted Child Vocalization Counts from LENA<sup>®</sup>’s output and contrasted them with clip-level counts of linguistic vocalizations by the key child in the human annotation.<sup>2</sup> As shown in Figure 4, there is a

---

<sup>2</sup>In a previous version of this analysis, we had calculated CVC as the number of CHN segments in LENA<sup>®</sup>, and the number of linguistic vocalizations as tagged by human annotators. Further inspection of LENA<sup>®</sup> documentation revealed this was incorrect, since LENA<sup>®</sup> counts can include several linguistic vocalizations within one CHN segment, and also includes linguistic vocalizations from CHF segments. Given

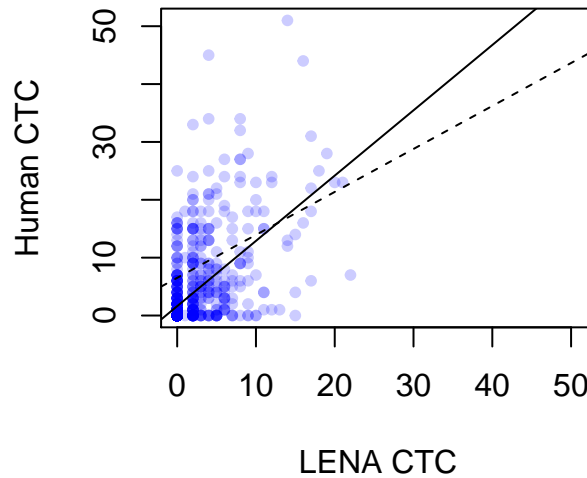
strong association between clip-level counts estimated via the LENA<sup>®</sup> system and those found in the human annotations: the Pearson correlation between the two was  $r = 0.73$  ( $p = < .001$ ) when all clips were taken into account, and  $r = 0.61$  ( $p = < .001$ ) when only clips with some child speech (i.e., excluding clips with 0 counts in both LENA<sup>®</sup> and human annotations) were considered (see Table 6 for correlation results of CVC, CTC, and AWC analyses). This suggests that the LENA<sup>®</sup> system captures differences in terms of number of child vocalizations across clips rather well.



*Figure 4.* Child Vocalization Counts according to LENA (x axis) and humans (y axis). Each point represents the CVC totaled within a clip. The solid line corresponds to a linear regression fit to data from all clips; the dashed line corresponds to an analysis excluding clips where both the human and LENA<sup>®</sup> found zero child vocalizations. The x and y ranges have been adjusted to be equal regardless of the data distribution.

The error analyses are reported on in Table 7. Generally speaking, LENA<sup>®</sup> tends to the inaccuracy of CHF, the latter decision seems potentially problematic. The same issue affected our CTC analyses. We now present analyses here that correctly represent LENA<sup>®</sup>'s reported CVC and CTC, since these are the field-standard measures. In Supplementary Materials (<https://osf.io/zdg6s>), we show results of the correlations and error analyses when CVC and CTC are calculated as the number of CHN/CHI segments instead. For CVC the results are identical; for CTC results were slightly worse results than those reported here.

underestimate vocalization counts, particularly when only clips with some speech are considered. This underestimation, however, is not systematic, and cumulating errors using the absolute error rate leads to an estimation of 73.77% deviation from the actual counts.



*Figure 5.* Conversational Turn Counts according to LENA (x axis) and humans (y axis). Each point represents the CTC totaled within a clip. The solid line corresponds to a linear regression fit to data from all clips; the dashed line corresponds to an analysis excluding clips where both the human and LENA<sup>®</sup> found zero child-adult or adult-child turns. The x and y ranges have been adjusted to be equal regardless of the data distribution for ease of visual comparison.

**Conversational Turn Counts (CTC) accuracy.** As with CVC, we extracted turn counts from LENA<sup>®</sup> logs, on the one hand, and attempted to do a similar calculation for the human annotations: the turn count is increased every time there is a vocalization by the key child followed by that of an adult (or vice versa) within 5 seconds. The association between clip-level LENA<sup>®</sup> and human CTC was weaker than that found for CVC (Figure 5): the Pearson correlation between the two was  $r = 0.57$  ( $p = < .001$ ) when all clips were taken into account, and  $r = 0.35$  ( $p = < .001$ ) when only clips with some child speech (i.e., excluding 208 clips with 0 counts in both LENA<sup>®</sup> and human annotations) were considered.

Error analyses can be found on Table 7. Inspection of errors and error rates reveals

Table 6

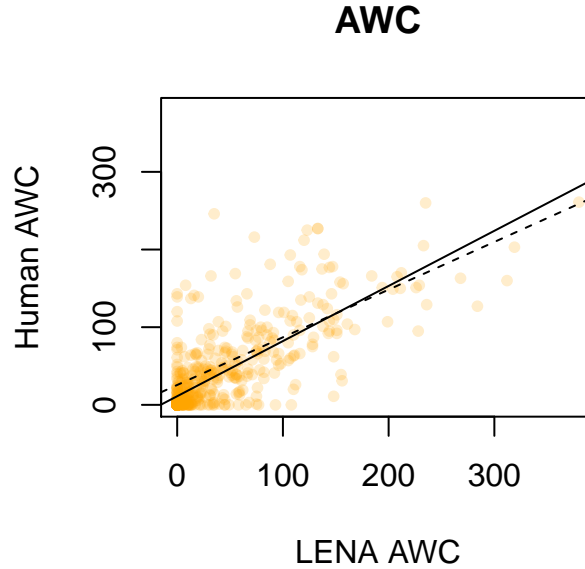
*Number of clips (N) and corresponding Pearson r coefficient for CVC, CTC, and AWC. 'N all' and 'r all' are computed over all clips. 'N' and 'r' represent non-null clips only (i.e., having some vocalizations, turns, and adult words respectively).*

	N all	r all	N	r
CVC	757	0.728	343	0.613
CTC	757	0.567	206	0.351
AWC	589	0.762	301	0.698

that LENA<sup>®</sup> tends to underestimate turn counts, which is particularly clear when excluding clips with no turns. As with vocalizations, the bias varied across clips leading to a cumulative absolute error rate of 93.61% deviation from the actual turn counts.

**Adult Word Counts accuracy.** One child in the (otherwise English) SOD corpus was learning French. Given our definition of orthographic words which is not language-specific, we have included this child to increase power, but results without them are nearly identical. See online Supplementary Materials, <https://osf.io/zdg6s>, for analyses excluding this child. In addition, a total of nine clips from three different WAR children contained some Spanish. Since we are uncertain of how accurate the transcriptions are for Spanish sentences, these clips were removed from consideration altogether.

The association between clip-level LENA<sup>®</sup> and human AWC in the four English-spoken corpora was strong (Figure 6): the Pearson correlation between the two was  $r=0.76$  ( $p<$



*Figure 6.* Adult Word Counts according to LENA (x axis) and humans (y axis). Each point represents the AWC totaled within a clip. The solid line corresponds to a linear regression fit to data from all clips; the dashed line corresponds to an analysis excluding clips where both the human and LENA<sup>®</sup> said there were no adult words. The x and y ranges have been adjusted to be equal regardless of the data distribution for ease of visual comparison.

.001) when all clips were taken into account, and  $r=0.70$  ( $p<.001$ ) when only clips with some adult speech (i.e., excluding 303 clips with 0 counts in both LENA<sup>®</sup> and human annotations) were considered. This suggests that the LENA<sup>®</sup> system captures differences in terms of number of adult word counts across clips well.

Table 7 contains error analyses for AWC, which reveal a different pattern from before. The error estimate across all clips was practically zero, with an error rate suggesting a tendency to over-estimate. Cumulating errors through the absolute error rate metric lead to an estimate of 123.76% deviation from the actual word counts.

**Effects of age and differences across corpora.** The preceding sections include overall results, with all but AWC conducted over all corpora. However, it is possible that performance would be higher for the corpora collected in North America (BER, WAR, SOD)



Table 7

*Mean (range) for each type of error estimate for CVC, CTC, and AWC. Error estimates are: E (error; NL-NH, where NL means the count according to LENA and NH the count according to the human), E-0 (error excluding clips with a zero count according to human or system analysis), ER (error rate;  $(NL-NH)/NH*100$ , in percent of the total), and AER (absolute ER;  $abs(NL-NH)/NH*100$ , in percent of the total, with abs meaning that we take the absolute); ER and AER exclude clips where the human count is zero.*

	E	(range)	E-0	(range)	ER %	(range)	AER %	(range)
CVC	-3	(-37,14)	-6	(-35,14)	-39	(-100,650)	74	(0,650)
CTC	-2	(-41,15)	-5	(-41,15)	-29	(-100,1200)	94	(0,1200)
AWC	-1	(-211,157)	-1	(-211,157)	54	(-100,7400)	124	(0,7400)

than those collected in other English-speaking countries (L05) or non-English speaking populations (TSI). Additionally, our age ranges are wide, and in the case of TSI children, some of the children are older than the oldest children in the LENA<sup>®</sup> training set. To assess whether accuracy varies as a function of corpora and child age, we fit mixed models. We report on key results here; for the full model output and additional analyses, please refer to our online Supplementary Materials (<https://osf.io/zdg6s>).

***Are there differences in false alarm, miss, and confusion rates as a function of corpus and child age?*** Figure 7 represents identification error rate as a function of age and corpus for individual children. A number of the children had a median identification error rate of zero due to the fact that they had many clips in which there was no speech, and LENA<sup>®</sup> had no false alarms, pulling the median to zero.

To test the possible impact of age and corpus statistically, we predicted false alarm, miss, and confusion rates in the analysis with all “Far” categories, TVN/ELE, and

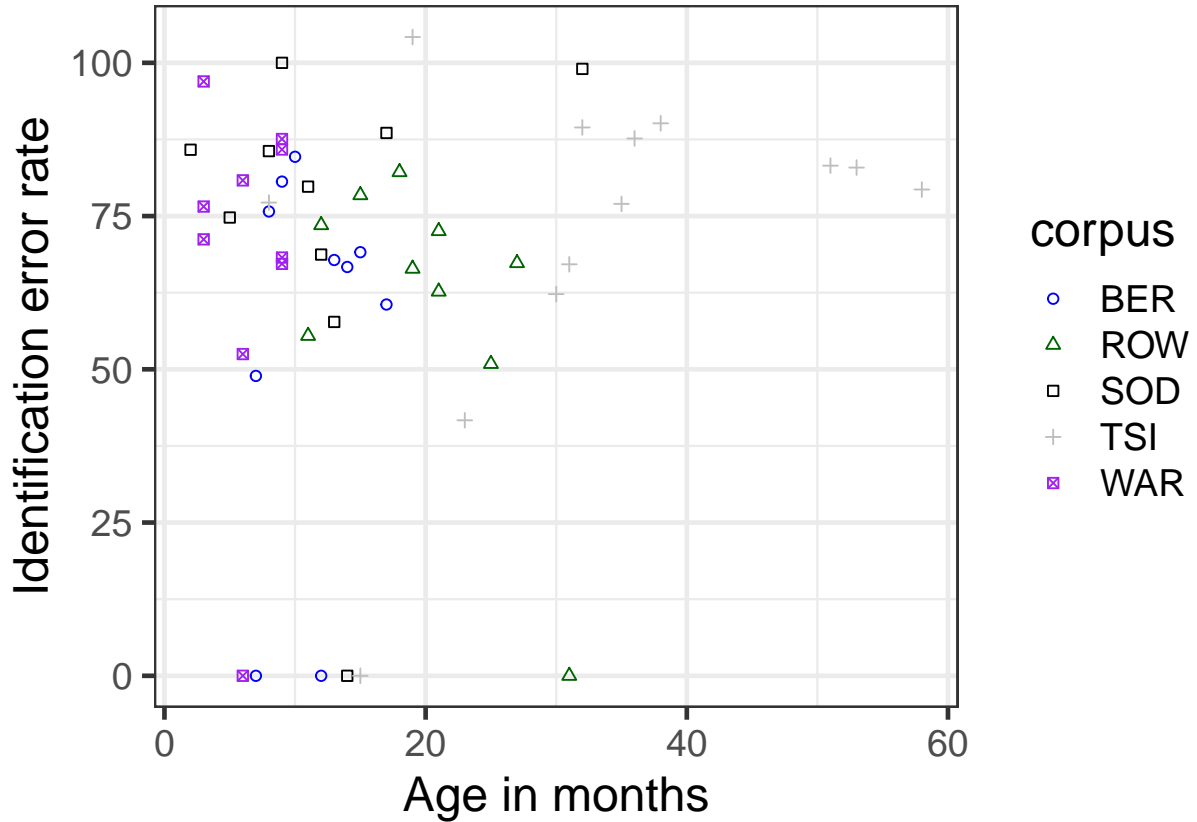


Figure 7. Identification error rate as a function of corpus and child age. Each point represents the median over all clips extracted from the data of one child. Color and shape indicates corpus: BER in blue circles, L05 is green triangles, SOD in black squares, TSI in gray pluses, and WAR in purple crossed squares.

OLN/OVL mapped onto Other (which yielded the best results in Section “False alarms, misses, confusion” above.) Our predictors were corpus, child age, and their interaction as fixed effects, and child ID as a random effect. We followed up with a Type III ANOVA to assess significance.

Across all analyses, corpus, child age, and their interaction were never significant, with the exception of confusion, where the interaction between corpus and age was significant at  $\alpha=.05$ . To investigate this effect further, we fit a mixed model predicting confusion rates from child age as fixed and child ID as random effects on each corpus separately. This

Table 8

*Results of Type III ANOVAs on false alarms (FA), misses (M), and confusions (C): Chi-square (degrees of freedom), followed by \* if the relevant factor is significant ( $p < .05$ ).*

	FA	M	C
Intercept	0.21 (1)	6.21 (1) *	4.88 (1) *
Corpus	1.98 (4)	3.16 (4)	3.65 (4)
Age	0.06 (1)	0.02 (1)	0.13 (1)
Corpus*Age	0.87 (4)	2.33 (4)	14.39 (4) *

revealed a main effect of age for SOD only (Chisq (1) = 14.53,  $p = < .001$ ; all other chi-squares were smaller than 14.53,  $p > .120$ ).

***Are there differences in CVC accuracy as a function of corpus and child age?*** For CVC, we fit a mixed model where manually-annotated CVC was predicted from LENA<sup>®</sup> CVC, in interaction with corpus and age, as fixed factors, and child ID as a random effect. Only effects and interactions involving the LENA<sup>®</sup> predictor are relevant to the present work. A Type III ANOVA found a two-way interaction suggesting that the predicted value of LENA<sup>®</sup> CVC with respect to manually-annotated CVC depended on the corpus. To investigate this further, we fit a model where manually-annotated CVC was predicted from LENA<sup>®</sup> CVC as a fixed effect and child ID as random factor, within each corpus separately. In the SOD corpus, we found a significant interaction between LENA<sup>®</sup> CVC and age (indicating that the predictive value of LENA<sup>®</sup> CVC decreased with child age); and a main effect of age (consistent with CVC going up with age in this corpus). The main effect of LENA<sup>®</sup> CVC, which is consistent with a predictive value of this metric with respect to manually-annotated CVC, emerged as a significant predictor in all corpora.

***Are there differences in CTC accuracy as a function of corpus and child age?*** For CTC, we fit a mixed model where CTC according to the human was predicted from CTC according to LENA<sup>®</sup>, in interaction with corpus and age, as fixed factors, declaring child ID as a random effect. This time our Type III ANOVA found a main effect of the LENA<sup>®</sup> CTC estimates, as well as an interaction between this factor and corpus. We followed up on this by fitting a model where CTC according to the human was predicted from CTC according to LENA<sup>®</sup> as fixed and child ID as random factor, for each corpus separately.

Inspection of these results (full output available from the Supplementary Materials, <https://osf.io/zdg6s>) suggests that the interaction emerged because the predictive value of LENA<sup>®</sup>'s CTC with respect to human counts was stronger for some corpora (Chi-squares for ROW 174.71, BER 107.49, and TSI 99.80) than others (Chi-squares for WAR 57.29, and SOD 30.91; all degrees of freedom are 1, and  $p < .001$ ).

***Are there differences in AWC accuracy as a function of corpus and child age?*** Finally, for AWC (which was only analyzable for the four ACLEW corpora), we fit a mixed model where AWC according to the human was predicted from AWC according to LENA<sup>®</sup>, in interaction with corpus and age, as fixed factors, declaring child ID as random effect. The Type III ANOVA revealed a three-way and both two-ways interactions involving the LENA<sup>®</sup> predictor, which was investigated by fitting additional mixed models to each corpus separately. An interaction between LENA<sup>®</sup> AWC and age was found for BER/WAR as well as SOD, due to a *decreased* predictive value of the LENA AWC with respect to the human AWC for older infants in BER and WAR but an *increase* in SOD. Notably, the positive association between LENA and human AWC was significant for all four corpora.

Table 9

*Results of Type III ANOVAs when predicting human counts (CVC, CTC, AWC) from LENA counts in interaction with age and corpus: Chi-square (degrees of freedom), followed by \* if the relevant factor is significant ( $p < .05$ ).*

	CVC	CTC	AWC
Intercept	5.34 (1) *	0.04 (1)	0 (1)
LENA	8.61 (1) *	19.1 (1) *	46.23 (1) *
Age	0.79 (1)	0.54 (1)	0.6 (1)
Corpus	8.18 (4)	3.41 (4)	1.71 (3)
LENA*Age	0 (1)	2.28 (1)	10.51 (1) *
LENA*Corpus	8.53 (4)	11.9 (4) *	15.7 (3) *
Age*Corpus	5.99 (4)	4.06 (4)	1.59 (3)
LENA*Age*Corpus	6 (4)	4.75 (4)	18 (3) *

## Discussion

The aim of the present study was to assess LENA<sup>®</sup> accuracy across key outcome measures: speaker classification accuracy, adult word counts, child vocalization counts, and conversational turn counts. We did this using an approach that sought to avoid inflating accuracy estimates in several ways. Methodologically, we used random or periodic sampling to select portions of the files for manual annotation, and our human annotators did not see the LENA<sup>®</sup> segmentation. Analytically, we considered both speech and non-speech classes (including electronic sounds and silence/Other). This permitted a systematic, extensive, and independent evaluation of LENA<sup>®</sup>'s key automated metrics. We also tested generalizability by analyzing LENA<sup>®</sup>'s performance across five different corpora: three based on the same population, language, dialect, and age group that LENA<sup>®</sup> was established for, and trained on

(North American English); one that allowed us to test how accurately it captured a different dialect of English (UK English); and one that tested its performance in a totally different recording situation (a rural setting with large families and many children present, speaking a linguistically unrelated language, and where the key children were, on average, somewhat older). We begin by recapping our key results.

Our first set of analyses tested overall accuracy, using established speech and talker segmentation metrics (false alarm rate, miss rate, confusion rate, and the composite identification error rate), and evaluated the pattern of errors in more detail, by assessing how LENA<sup>®</sup> and human annotators agreed (precision and recall). The identification error rate was relatively high (global 73%), mainly due to a high miss rate (missing or excluding speech that was there; 39%). The false alarm rate (identifying non-speech/silence as speech; 6%) and confusion rate (identifying voice type; 8%) were low.

To put these numbers in context, we asked the ACLEW project members to share with us preliminary results of their inter-rater reliability study. This study covers six corpora, including the four ACLEW corpora used here. For the present analysis, they considered the “gold” to be the original complete annotations, and the “system” the reliability annotations, which were done later and in only a subset of the corpus (one minute per child, for a total of 60 minutes across their six corpora). While we cannot report on these results in full because their publication is intended elsewhere, we can state the following overall observations, that can be compared against the LENA data reported on below. Among two human annotators, the ACLEW team reported an identification error rate of 56% (due to 20% false alarms, 19% miss rates, and 17% of confusion); for the four databases included here, the average identification error rate was 47%. This is considerably lower than the identification error rates reported for LENA<sup>®</sup> here, mainly due to a much lower miss rates, whereas both false alarm rates and confusion rates are higher across the two human coders. Inspection of false alarms and misses suggests the disagreement across humans emerges when there is

background speech, that one coder may pick up on and not the other.<sup>3</sup>

Another question is how this fares compared to other automatic systems. Our thorough review of the literature revealed that no previous report is comparable: Most often, the data used is considerably different (and overall easier; e.g., recorded in formal settings, with a small number of speakers, who produce long vocalizations); moreover, previous research tends to overestimate performance by using lax evaluation criteria (e.g., allowing errors in a “collar” around each vocalization). The most comparable data point comes from the DIHARD Challenge (Ryant et al., 2019). DIHARD employed data from a range of domains, including daylong recordings; in fact, they used a different selection of data from the BER corpus used here. The subset of BER used for DIHARD is likely to lead to lower error rates because they selected only files that contained some speech; by excluding files with little to no speech, they prevent the appearance of very high diarization error rates (which emerge when the numerator, i.e. the amount of speech, is very small). Thus, the DIHARD reanalyses are likely to overestimate the systems’ performance in terms of data selection. Their evaluation, however, was comparable to the one we used here, with no leeway or collar. Diarization error rates for the BER subset by systems submitted to DIHARD 2019 varied between 48% and 121%, with a median around 70%. Thus, LENA<sup>®</sup> is competitive with respect to state-of-the-art systems, although some of them do score considerably better.<sup>4</sup>

Returning to the LENA<sup>®</sup> system results, the overall error rate can be fruitfully interpreted by considering performance on individual speaker tags. In terms of precision (to what extent do LENA<sup>®</sup> tags contain what they say they contain), the system performed

---

<sup>3</sup>Taking all categories together, Cohen’s  $\kappa$  agreement was .64 (weighted  $\kappa$  .65) for the ACLEW inter-rater reliability coding on all six ACLEW datasets, which is higher than the best case scenario for LENA.

<sup>4</sup>DIHARD uses diarization error rate on individual speakers’ identities, rather than identification error rates on speaker types as we do here. There is no mathematical procedure to derive one from the other, except in the case when there is one speaker per speaker type, in which case diarization error rate is most likely identical to identification error rate.

relatively well at identifying female voices (60% of frames tagged by LENA<sup>®</sup> as FAN were coded as female adult by the human coders), and the target child (61% of frames tagged by LENA<sup>®</sup> as CHN were correct). However, the system performed substantially worse with other talker types (e.g. 43% and 28% for MAN and CXN, respectively); that is, less than a half of the frames that LENA<sup>®</sup> tagged as being speech spoken by these speakers actually correspond to them.

To get a sense of how these results compare to multiple human coders, we also asked about precision and recall in the reliability data from the ACLEW team. Across all six corpora, precision for key child was the highest, at 80%; for the other speakers it was: 72% female adult, 72% male adult, and 65% other child. Precision is higher and more similar across speaker types in the ACLEW reliability data than in our LENA<sup>®</sup>-human comparison here.

In terms of recall (how accurately LENA<sup>®</sup> captured the human annotations), performance for the key child’s vocalizations was moderately robust: 50% of the frames humans attributed to the key child speech were captured by LENA<sup>®</sup> under the CHN tag. However, recall was poorer for the other three talker types, at around 31-32%. As for recall in the ACLEW reliability data, the key child score was 79%; for the other speakers it was: 71% female adult, 63% male adult, and 55% other child. Thus, although we see lower recalls for male adults and other children in both, the overall level of recall is much higher across two human coders than between LENA<sup>®</sup> and human, mainly due to LENA<sup>®</sup>’s tendency to miss speech more than humans do.

Our second set of analyses tested the accuracy of three of the aggregated counts automatically provided by LENA<sup>®</sup>, namely Child Vocalization Counts (CVC), Conversational Turn Counts (CTC) and Adult Word Counts (AWC). We found relatively high correlations between clip-level counts estimated via the LENA<sup>®</sup> system and those from the human annotations for AWC and CVC, with weaker performance for CTC.



However, such correlational analyses do not establish whether LENA<sup>®</sup> systematically over- or under-estimates. For this we examined several types of error estimates. For overall error estimates (i.e., how far was the LENA<sup>®</sup> count from the human annotators'), the means across clips for CVC, CTC, and AWC was an encouraging -3.32, -1.85, and -1.04, respectively. These low errors were not solely due to many clips lacking vocalizations, turns, or adult words altogether, because when we exclude such clips we still get what seem to be low errors: means were -6.46, -5.08, and -0.67 for CVC, CTC, and AWC respectively.

We also put these error patterns in context by taking into account how large the counts were to begin with. Such error rates, however, are only defined for files which, according to the human, contain at least one unit (otherwise, we divide an error of a certain size by zero, which is undefined). We find error rates suggesting that LENA<sup>®</sup> counts are off by between a third or a half of the original counts. Inspection of the sign in these rates indicates that, by and large, LENA<sup>®</sup> systematically underestimates the raw counts of its main quantitative measures - particularly child vocalizations and conversational turns, and to a lesser extent, adult words, which showed more erratic error patterns. In addition, we find much larger absolute error rates (i.e. in our analysis that prevents under- and over-estimations from cancelling each other out.)

Finally, we also inspected the extent to which LENA<sup>®</sup> performance was affected by dialect, language, and child age in a final set of analyses. We would like to be tentative about the interpretation of these results, because we only have about 10 children, often varying widely in age, in each corpus, with some mismatch in age range across corpora (see Table 2). This means that we did not have a great deal of power to capture true differences across corpora and that we may have some spurious effects or interactions due to chance differences.

With all these caveats in mind, we predicted that performance would be higher for the corpora collected in North America (BER, WAR, SOD) than for corpora collected in other English-speaking countries (L05) or non-English speaking populations (TSI), and that

accuracy would decrease with age, since our sample contains children older than those included in the LENA<sup>®</sup> training set. This is not what we found. For instance, we found an interaction between corpus and age for the confusion rate, due to an increase in confusion errors within the SOD corpus but not in any of the others, a result that we have no ready explanation for. Similarly, LENA<sup>®</sup> counts predicted human counts in the CVC, CTC, and AWC analyses, and although we did observe some interactions, none of them were easy to interpret and none explained away the predictive value. As just mentioned, we are cautious in this interpretation, and invite further work on bigger samples (more data per child, more children per corpus) to ensure sufficient power and precision.

In general, whether LENA<sup>®</sup> results are “good enough” for a given research, educational, or clinical study depends largely on the goals of each particular study. For example, we can describe precision rates of 60% (i.e., 60% of frames tagged by LENA<sup>®</sup> as FAN were coded as female adult by human coders) and 61% (i.e., 61% of frames tagged as target child were also tagged as such by human coders) as being reasonably good. We do this partly because these rates are much higher than the system’s precision rates for other speakers (MAN 43%, CXN 28%) but also partly because our frame-based criteria is more stringent than many coding schemes previously applied. That said, whether a *particular* accuracy rate can be considered sufficient will depend on the purpose of the study. As a result, we next provide a set of recommendations to help researchers make this determination for their goals.

**What research goals can one pursue given the performance of LENA<sup>®</sup> segmentation and metrics?** In the present corpora, the system’s false alarm rate (i.e., identifying speech where there was none) was very low while its miss rate (missing speech that was actually there) was relatively high. This makes LENA<sup>®</sup> more suitable for studies in which it is extremely important not to “invent” speech that is not there but less suitable for studies in which capturing most, if not all, of the speech produced is crucial. Based on these findings, LENA<sup>®</sup> would be a good tool for finding “high talk volume” parts of the day for a)

careful further transcription (e.g. of low-frequency events like a certain grammatical construction of interest), b) annotation of specific speech characteristics (e.g. mean length of utterance), or c) comparing relative talk volume across samples. However, we advise caution in using LENA<sup>®</sup> when raw quantity of speech is crucial for the research question, or when small differences in talk volume might have very significant theoretical consequences; this is often the case in clinical populations where children’s own vocalizations can be an important diagnosis-relevant characteristic (e.g., in children who are deaf or hard of hearing, individuals with ASD, speech apraxia, etc.).

Similarly, although the overall confusion rate (i.e. incorrectly identifying talkers, such as giving a “female adult” tag for a “child” utterance) for LENA<sup>®</sup> was very low, this does not fully convey the level of accuracy for speech, particularly when considering every talker type. In terms of precision, the system’s female adult and key child categorization was quite accurate, whereas precision was lower for male adults and other children: the vast majority of the frames labeled as male adult or other children did not in fact contain speech by these speaker types. In terms of recall, LENA<sup>®</sup> was good at capturing speech by the key child as such, but recall was lower for the other talker categories.

We, thus, recommend caution before undertaking any analyses that rely on the accuracy (precision and/or recall) of male adult and other children’s speech. For example, if the goal is simply to calculate an overall adult word count (AWC), summing over male and female adult speakers, some confusion between MAN and FAN is likely not problematic. However, if the goal of the study is to compare the relative input from fathers and mothers, LENA<sup>®</sup> tags are relatively unreliable and on our view, merit further manual vetting in most use cases.

As another example (detailed further in the “Recall” results above), if the goal is to capture as many of the key child’s vocalisations as possible, it might be worthwhile to pull out segments LENA<sup>®</sup> labelled as non-target child, CXN, (of which 23% was target child

speech) as well, with human coders brought in to filter out non-target child speech. Indeed, we find that this kind of binary classification (key child or not) can be readily undertaken with little training by research assistants in our labs, and would substantially boost data quality and quantity for child vocalizations in this use case.

Notably, while we recommend LENA<sup>®</sup> users be cautious in their use of LENA<sup>®</sup> identification and classification, especially for certain talker classes, our results for LENA<sup>®</sup> count metrics suggest these derived counts may be accurate enough to serve well across a large variety of uses. To begin with, as far as it is possible to generalize from the limited range of samples tested here (children aged 2 to 58 months, learning North American English, UK English, or Tsimane’) it seems that the system’s performance does not vary a great deal across ages, dialects, language and home settings. Moreover, *correlations* between human and LENA<sup>®</sup> clip-level counts were high to very high, suggesting that the software accurately captures differences in counts across clips (even when *error* rates were also high). These correlations remained quite high even when clips with counts equal to zero were removed from consideration, suggesting that LENA<sup>®</sup> captures gradience in vocalization counts.

However, our finding that LENA<sup>®</sup> generally underestimates the quantity of child vocalizations, child-adult turns, and adult words deserves further consideration. Indeed, further work is needed to fully understand the nature and extent of this limitation. Our clips were 1–2 minutes in length, and therefore they either tended to have very little speech or a lot of it. Error rates over hours could be smaller, because local errors average out; or greater, if the LENA<sup>®</sup> system systematically underestimates counts. In a LENA<sup>®</sup> technical report, AWC accuracy was variable across two 12-hour recordings: 1% lower than human transcription for one child, but 27% lower for a second child. This same report notes that AWC accuracy quickly plateaus as recording time increases beyond one hour, leveling to 5-10% in recordings greater than 2 hours in length (Xu et al., 2009).

Thus, it is important for further work to help establish the systematicity in the

estimates provided by LENA<sup>®</sup>: if underestimates are robust and systematic (as suggested by present results for CVC and CTC, but not AWC), it may be possible to develop a correction factor to compensate for this bias. However, this bias may be challenging to nail down precisely and for AWC it may depend on the language in question. For instance, a recent study in Finland documented that LENA<sup>®</sup> largely *overestimated* AWC and only slightly underestimated child vocalization counts (Elo, 2016).

### **How to test the reliability of the automated output provided by LENA<sup>®</sup>.**

We are overall hopeful that the reliability metrics we provide here will be relevant for researchers working with different populations. We hope the current paper inspires others to evaluate and report all aspects of the system, rather than a subset of metrics. Similarly extensive evaluations of LENA<sup>®</sup> in other corpora would bolster the validation literature, and be useful for the whole research community. In fact, it would be useful if researchers systematically test the reliability of LENA<sup>®</sup> counts in their own samples, especially if they are collecting data from families living in different environments from those assessed here. Next, we provide some guidelines for how to go about this. Note that this requires downloading the audio (.wav) file generated by LENA<sup>®</sup> as well as the corresponding LENA<sup>®</sup> output file.

First, we recommend a literature search, to determine whether a similar sample has been studied in the past for which there exists reliability data (see for example, Cristia et al., 2019 for a systematic review). If no studies exist, draw 10 x 2 minutes randomly from 10 children. This is about 3h20min of data, which takes roughly 60h to annotate, in our experience. We recommend training annotators using the ACLEW Annotation Scheme <https://osf.io/b2jep/>, which has an online test annotators can go through to ensure reliability. Once the manual annotations are complete, the LENA<sup>®</sup> annotations can be extracted and compared against the human annotation using the free DiViMe software ([divime.readthedocs.io](http://divime.readthedocs.io), Le Franc et al., 2018). This will allow researchers to extract the

classification accuracy measures used here (false alarm rate, miss rate, confusion rate and the derived identification error rate), as well as CVC, CTC, and AWC comparing LENA<sup>®</sup> and human annotations. We note this DiViMe pipeline is only possible “off the shelf” for manual annotations made using the ACLEW Annotation Scheme, though in principle, it is adaptable to other schemata by adept programmers.

One issue that may arise is whether data should be sampled differently to, for example, make sure every class is represented the same amount of time and/or a minimum of time. Our understanding is that class imbalance and data scarceness is an important issue for training, and directly affects algorithm accuracy (this is a general problem, but to cite just one example on GMMs, Garcia-Moral, Solera-Urena, Pelaez-Moreno, & Diaz-de-Maria, 2011). However, it does not pose the same kind of problem for evaluation. That is, if there are no samples of a given category, then accuracy cannot be evaluated; if there are only a few, then it is possible that these are special in some way and accuracy estimates may not generalize well to others. Thus, it would indeed be desirable to have enough samples of a given label to reduce the impact of each individual instance, in case they are outliers. That said, almost any strategy that attempts to boost the frequency of specific categories risks increasing the issue of non-generalizability. For instance, if one were to over-sample regions tagged by LENA<sup>®</sup> as MAN in the hopes of having more male samples, one may only be capturing certain types of male speech or acoustic properties. To take this example further, notice that male speech is our smallest category, representing 1% of the data. Since we sampled randomly or periodically, this represents the prevalence of male speech and the samples that are included are unlikely to be acoustically biased.

Separately, researchers should reflect about the accuracy needed for their question of interest. For instance, suppose we have an evaluation of an intervention where we expect treatment children to hear 20% more speech than controls, or an individual difference study where we expect that the lower fifth of the children hear 20% less speech than the top fifth.

If the intended measure used to compare groups has an error rate larger than the effect predicted (such as the the CTC error rate we find here), a different algorithm or outcome metric would be wise.

## Conclusions

In conclusion, in this study, we have provided a broad evaluation of accuracy across the key outcome measures provided by LENA<sup>®</sup> (classification, child vocalization counts, conversational turn counts, and adult word counts), and its generalizability across different dialects, languages, ages, and settings. We have provided some recommendations for how to use LENA<sup>®</sup> in future studies most effectively, and how to test the accuracy of the LENA<sup>®</sup> algorithms on particular samples of data.

There are, however, a number of areas of research that we have not addressed. For example, we have not investigated how accurately LENA<sup>®</sup> detects individual variation across children or families. It would be particularly useful to know whether LENA<sup>®</sup> can classify children with the sensitivity and specificity needed for accurate identification of language disorders. Oller et al. (2010) used LENA<sup>®</sup> to differentiate vocalizations from 232 typically developing children and children with autism or language delay with a high degree of accuracy. However, key to this was the use of additional algorithms, not yet available from LENA<sup>®</sup>, to identify and classify the acoustic features of “speech-related vocal islands”. Further work (including shared code) would greatly bolster progress on this topic.

Even if it turns out that LENA<sup>®</sup> is not accurate enough to classify children precisely for a given ability or diagnosis, it may be accurate enough to capture the rank order of individual children’s language growth, which can provide useful information about the relative language level of children in a sample or population (see, e.g., Gilkerson et al., 2017). Similarly, LENA<sup>®</sup> may not accurately capture the precise number of child vocalisations

produced over time, but it may track developmental trajectory (e.g., the slope of growth) relatively well. Finally, although our results suggest that aspects of the system’s output may be relatively robust to differences across languages and dialects, we need more evidence of how it fares across mono- and multi-lingual language environments (cf. Orena, 2019).

It is undeniable that children learn language from the world around them. Naturalistic daylong recordings offer an important avenue to examine this uniquely human development, alongside other fundamental questions about human interaction, linguistic typology, psychology, and sociology. Tools and approaches that allow us to tap such recordings’ contents stand to contribute deeply to our understanding of these processes. We look forward to further work that addresses the many remaining questions within this area.

## Acknowledgments

This research benefits from the Analyzing Child Language Experiences around the World (ACLEW) collaborative project funded by the Trans-Atlantic Platform for Social Sciences and Humanities “Digging into Data” challenge, including a local Academy of Finland grant (312105) to OR, ANR-16-DATA-0004 ACLEW to AC, NEH HJ-253479-17 to EB, and funding from the Social Sciences and Humanities Research Council of Canada (869-2016-0003) and the Natural Sciences and Engineering Research Council of Canada (501769-2016-RGPDD) to MS. AC acknowledges further support from (ANR-17-CE28-0007 LangAge, ANR-14-CE30-0003 MechELex, ANR-17-EURE-0017); and the James S. McDonnell Foundation Understanding Human Cognition Scholar Award. MS was also funded by a Social Sciences and Humanities Research Council of Canada Insight Grant (435-2015-0628). EB acknowledges NIH (DP5 OD019812-01). CR was also funded by the Economic and Social Sciences Research Council (ES/L008955/1). OR was also funded by an Academy of Finland grant no. 314602.



## Open Practices Statement

The study relies indirectly on daylong audiorecordings (which cannot be made public to protect participants) and human and LENA<sup>®</sup> annotations for extracted clips (which are not deidentified); these are stored in private repositories that do not have a persistent identifier. The annotation data were used to generate statistics at the clip level, which are the input to analyses presented here. Both the clip level statistics and all analyses are publicly available from [https://github.com/jsalt-coml/lena\\_eval](https://github.com/jsalt-coml/lena_eval). None of these analyses were pre-registered.

## References

- Baumer, B., Cetinkaya-Rundel, M., Bray, A., Loi, L., & Horton, N. J. (2014). R Markdown: Integrating a reproducible analysis tool into introductory statistics. *arXiv Preprint arXiv:1402.1894*.
- Bergelson, E. (2016). Bergelson Seedlings Homebank corpus. <https://doi.org/10/T5PK6D>
- Bergelson, E., Casillas, M., Soderstrom, M., Seidl, A., Warlaumont, A. S., & Amatuni, A. (2019). What do North American babies hear? A large-scale cross-corpus analysis. *Developmental Science*, 22(1), e12724.
- Bergelson, E., Cristia, A., Soderstrom, M., Warlaumont, A., Rosenberg, C., Casillas, M., ... Bunce, J. (2017). ACLEW project. Databrary.
- Bredin, H. (2017). Pyannote.metrics: A toolkit for reproducible evaluation, diagnostic, and error analysis of speaker diarization systems. In *INTERSPEECH* (pp. 3587–3591).
- Bulgarelli, F., & Bergelson, E. (2019). Look who’s talking: A comparison of automated and human-generated speaker tags in naturalistic day-long recordings. *Behavior Research Methods*, 1–13.
- Busch, T., Sangen, A., Vanpoucke, F., & Wieringen, A. van. (2018). Correlation and agreement between Language ENvironment Analysis (LENATM) and manual transcription for Dutch natural language recordings. *Behavior Research Methods*, 50(5), 1921–1932. <https://doi.org/10.3758/s13428-017-0960-0>
- Canault, M., Le Normand, M. T., Foudil, S., Loundon, N., & Thai-Van, H. (2016). Reliability of the Language ENvironment Analysis system (LENATM) in European French. *Behavior Research Methods*, 48(3), 1109–1124. <https://doi.org/10.3758/s13428-015-0634-8>

- Casillas, M., Bergelson, E., Warlaumont, A. S., Cristia, A., Soderstrom, M., VanDam, M., & Sloetjes, H. (2017). A new workflow for semi-automatized annotations: Tests with long-form naturalistic recordings of children's language environments. In *Interspeech 2017* (pp. 2098–2102).
- Cristia, A., Bulgarelli, F., & Bergelson, E. (2019). Accuracy of the Language Environment Analysis System: A systematic review. Retrieved from <https://osf.io/fhs57>
- d'Apice, K., Latham, R. M., & Stumm, S. von. (2019). A naturalistic home observational approach to children's language, cognition, and behavior. *Developmental Psychology*.
- Elo, H. (2016). *Acquiring language as a twin*. Tampere, Finland: Tampere University Press.
- Ganek, H. V., & Eriks-Brophy, A. (2018). A concise protocol for the validation of Language ENvironment Analysis (LENA) conversational turn counts in vietnamese. *Communication Disorders Quarterly*, 39(2), 371–380.
- Garcia-Moral, A. I., Solera-Urena, R., Pelaez-Moreno, C., & Diaz-de-Maria, F. (2011). Data Balancing for Efficient Training of Hybrid ANN/HMM Automatic Speech Recognition Systems. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(3), 468–481. <https://doi.org/10.1109/TASL.2010.2050513>
- Gilkerson, J., Coulter, K. K., & Richards, J. A. (2008). Transcriptional analyses of the LENA natural language corpus. LENA Foundation.
- Gilkerson, J., & Richards, J. A. (2008). The LENA Natural Language Study. LENA Foundation.
- Gilkerson, J., Richards, J. A., Warren, S. F., Montgomery, J. K., Greenwood, C. R., Kimbrough Oller, D., ... Paul, T. D. (2017). Mapping the early language environment using all-day recordings and automated analysis. *American Journal of*

- Speech-Language Pathology*, 26(2), 248.  
[https://doi.org/10.1044/2016\\_AJSLP-15-0169](https://doi.org/10.1044/2016_AJSLP-15-0169)
- Gilkerson, J., Zhang, Y., Xu, D., Richards, J. A., Xu, X., Jiang, F., . . . Toppings, K. (2016). Evaluating language environment analysis system performance for Chinese: A pilot study in Shanghai. *Journal of Speech Language and Hearing Research*, 85(2), 445–452. <https://doi.org/10.1044/2015>
- Goh, K.-I., & Barabási, A.-L. (2008). Burstiness and memory in complex systems. *EPL (Europhysics Letters)*, 81(4), 48002.
- Greenwood, C. R., Thiemann-Bourque, K., Walker, D., Buzhardt, J., & Gilkerson, J. (2011). Assessing children’s home language environments using automatic speech recognition technology. *Communication Disorders Quarterly*, 32(2), 83–92.  
<https://doi.org/10.1177/1525740110367826>
- Lamere, P., Kwok, P., Gouvea, E., Raj, B., Singh, R., Walker, W., . . . Wolf, P. (2003). The CMU SPHINX-4 speech recognition system. In *IEEE Intl. Conf. on Acoustics, Speech and Signal Processing*. Hong Kong.
- Le Franc, A., Riebling, E., Karadayi, J., Wang, Y., Scaff, C., Metze, F., . . . others. (2018). The aclew divime: An easy-to-use diarization tool. In *Interspeech* (pp. 1383–1387).
- Lehet, M., Arjmandi, M. K., Dilley, L. C., Roy, S., & Houston, D. (2018). Fidelity of automatic speech processing for adult speech classifications using the Language ENvironment Analysis (LENA) system. *Proceedings of Interspeech*, 3–7.
- MacWhinney, B. (2017). Tools for Analyzing Talk Part 1: The CHAT Transcription Format. Carnegie.
- McDivitt, K., & Soderstrom, M. (2016). McDivitt homebank corpus.

- Oller, D. K., Niyogi, P., Gray, S., Richards, J. A., Gilkerson, J., Xu, D., . . . Cutler, E. A. (2010). Automated vocal analysis of naturalistic recordings from children with autism, language delay, and typical development. *Proceedings of the National Academy of Sciences*, 107(30), 13354–13359. <https://doi.org/10.1073/pnas.1003882107>
- Orena, A. J. (2019). Growing up bilingual: Examining the language input and word segmentation abilities of bilingual infants. PsyArXiv. <https://doi.org/10.31234/osf.io/x9wr8>
- Rowland, C. F., Bidgood, A., Durrant, S., Peter, M., & Pine, J. M. (2018). The Language 0-5 Project. University of Liverpool. <https://doi.org/10.17605/OSF.IO/KAU5F>
- RStudio Team. (2019). *RStudio: Integrated development environment for R*. Boston, MA: RStudio, Inc. Retrieved from <http://www.rstudio.com/>
- Ryant, N., Church, K., Cieri, C., Cristia, A., Du, J., Ganapathy, S., & Liberman, M. (2019). Second DIHARD Challenge evaluation plan. *Linguistic Data Consortium, Tech. Rep.*
- Scaff, C., Stieglitz, J., Casillas, M., & Cristia, A. (n.d.). Daylong audio recordings of young children in a forager-farmer society show low levels of verbal input with minimal age-related change.
- Seidl, A., Cristia, A., Soderstrom, M., Ko, E.-S., Abel, E. A., Kellerman, A., & Schwichtenberg, A. (2018). Infant-mother acoustic-prosodic alignment and developmental risk. *Journal of Speech, Language, and Hearing Research*, 61(6), 1369–1380.
- Soderstrom, M., Bergelson, E., Warlaumont, A., Rosemberg, C., Casillas, M., Rowland, C., . . . Bunce, J. (n.d.). The ACLEW Random Sampling corpus.
- Team, R. C., & others. (2013). R: A language and environment for statistical computing.

Vienna, Austria.

- VanDam, M., & De Palma, P. (2018). A modular, extensible approach to massive ecologically valid behavioral data. *Behavior Research Methods*.  
<https://doi.org/10.3758/s13428-018-1167-8>
- VanDam, M., Warlaumont, A. S., Bergelson, E., Cristia, A., Soderstrom, M., De Palma, P., & MacWhinney, B. (2016). HomeBank: An online repository of daylong child-centered audio recordings. In *Seminars in speech and language* (Vol. 37, pp. 128–142). Thieme Medical Publishers.
- Warlaumont, A., Pretzer, G., Walle, E., Mendoza, S., & Lopez, L. (2016). Warlaumont HomeBank corpus.
- Weisleder, A., & Fernald, A. (2013). Talking to children matters. *Psychological Science*, 24(11), 2143–2152. <https://doi.org/10.1177/0956797613488145>
- Xu, D., Yapanel, U., & Gray, S. (2009). Reliability of the LENATM Language Environment Analysis System in young children’s natural home environment. LENA Foundation.
- Zimmerman, F. J., Gilkerson, J., Richards, J. A., Christakis, D. A., Xu, D., Gray, S., & Yapanel, U. (2009). Teaching by listening: The importance of adult-child conversations to language development. *Pediatrics*, 124(1), 342–349.
- Baumer, B., Cetinkaya-Rundel, M., Bray, A., Loi, L., & Horton, N. J. (2014). R Markdown: Integrating a reproducible analysis tool into introductory statistics. *arXiv Preprint arXiv:1402.1894*.
- Bergelson, E. (2016). Bergelson Seedlings Homebank corpus. <https://doi.org/10/T5PK6D>
- Bergelson, E., Casillas, M., Soderstrom, M., Seidl, A., Warlaumont, A. S., & Amatuni, A. (2019). What do North American babies hear? A large-scale cross-corpus analysis.

*Developmental Science*, 22(1), e12724.

- Bergelson, E., Cristia, A., Soderstrom, M., Warlaumont, A., Rosemberg, C., Casillas, M., ... Bunce, J. (2017). ACLEW project. Databrary.
- Bredin, H. (2017). Pyannote.metrics: A toolkit for reproducible evaluation, diagnostic, and error analysis of speaker diarization systems. In *INTERSPEECH* (pp. 3587–3591).
- Bulgarelli, F., & Bergelson, E. (2019). Look who's talking: A comparison of automated and human-generated speaker tags in naturalistic day-long recordings. *Behavior Research Methods*, 1–13.
- Busch, T., Sangen, A., Vanpoucke, F., & Wieringen, A. van. (2018). Correlation and agreement between Language ENvironment Analysis (LENATM) and manual transcription for Dutch natural language recordings. *Behavior Research Methods*, 50(5), 1921–1932. <https://doi.org/10.3758/s13428-017-0960-0>
- Canault, M., Le Normand, M. T., Foudil, S., Loundon, N., & Thai-Van, H. (2016). Reliability of the Language ENvironment Analysis system (LENATM) in European French. *Behavior Research Methods*, 48(3), 1109–1124. <https://doi.org/10.3758/s13428-015-0634-8>
- Casillas, M., Bergelson, E., Warlaumont, A. S., Cristia, A., Soderstrom, M., VanDam, M., & Sloetjes, H. (2017). A new workflow for semi-automatized annotations: Tests with long-form naturalistic recordings of children's language environments. In *Interspeech 2017* (pp. 2098–2102).
- Cristia, A., Bulgarelli, F., & Bergelson, E. (2019). Accuracy of the Language Environment Analysis System: A systematic review. Retrieved from <https://osf.io/fhs57>
- d'Apice, K., Latham, R. M., & Stumm, S. von. (2019). A naturalistic home observational

- approach to children's language, cognition, and behavior. *Developmental Psychology*.
- Elo, H. (2016). *Acquiring language as a twin*. Tampere, Finland: Tampere University Press.
- Ganek, H. V., & Eriks-Brophy, A. (2018). A concise protocol for the validation of Language ENvironment Analysis (LENA) conversational turn counts in vietnamese. *Communication Disorders Quarterly*, 39(2), 371–380.
- Garcia-Moral, A. I., Solera-Urena, R., Pelaez-Moreno, C., & Diaz-de-Maria, F. (2011). Data Balancing for Efficient Training of Hybrid ANN/HMM Automatic Speech Recognition Systems. *EEE Transactions on Audio, Speech, and Language Processing*, 19(3), 468–481. <https://doi.org/10.1109/TASL.2010.2050513>
- Gilkerson, J., Coulter, K. K., & Richards, J. A. (2008). Transcriptional analyses of the LENA natural language corpus. LENA Foundation.
- Gilkerson, J., & Richards, J. A. (2008). The LENA Natural Language Study. LENA Foundation.
- Gilkerson, J., Richards, J. A., Warren, S. F., Montgomery, J. K., Greenwood, C. R., Kimbrough Oller, D., ... Paul, T. D. (2017). Mapping the early language environment using all-day recordings and automated analysis. *American Journal of Speech-Language Pathology*, 26(2), 248. [https://doi.org/10.1044/2016\\_AJSLP-15-0169](https://doi.org/10.1044/2016_AJSLP-15-0169)
- Gilkerson, J., Zhang, Y., Xu, D., Richards, J. A., Xu, X., Jiang, F., ... Toppings, K. (2016). Evaluating language environment analysis system performance for Chinese: A pilot study in Shanghai. *Journal of Speech Language and Hearing Research*, 85(2), 445–452. <https://doi.org/10.1044/2015>
- Goh, K.-I., & Barabási, A.-L. (2008). Burstiness and memory in complex systems. *EPL*



- (*Europhysics Letters*), 81(4), 48002.
- Greenwood, C. R., Thiemann-Bourque, K., Walker, D., Buzhardt, J., & Gilkerson, J. (2011). Assessing children's home language environments using automatic speech recognition technology. *Communication Disorders Quarterly*, 32(2), 83–92.  
<https://doi.org/10.1177/1525740110367826>
- Lamere, P., Kwok, P., Gouvea, E., Raj, B., Singh, R., Walker, W., ... Wolf, P. (2003). The CMU SPHINX-4 speech recognition system. In *IEEE Intl. Conf. on Acoustics, Speech and Signal Processing*. Hong Kong.
- Le Franc, A., Riebling, E., Karadayi, J., Wang, Y., Scaff, C., Metze, F., ... others. (2018). The aclew divime: An easy-to-use diarization tool. In *Interspeech* (pp. 1383–1387).
- Lehet, M., Arjmandi, M. K., Dilley, L. C., Roy, S., & Houston, D. (2018). Fidelity of automatic speech processing for adult speech classifications using the Language ENvironment Analysis (LENA) system. *Proceedings of Interspeech*, 3–7.
- MacWhinney, B. (2017). Tools for Analyzing Talk Part 1: The CHAT Transcription Format. Carnegie.
- McDivitt, K., & Soderstrom, M. (2016). McDivitt homebank corpus.
- Oller, D. K., Niyogi, P., Gray, S., Richards, J. A., Gilkerson, J., Xu, D., ... Cutler, E. A. (2010). Automated vocal analysis of naturalistic recordings from children with autism, language delay, and typical development. *Proceedings of the National Academy of Sciences*, 107(30), 13354–13359. <https://doi.org/10.1073/pnas.1003882107>
- Orena, A. J. (2019). Growing up bilingual: Examining the language input and word segmentation abilities of bilingual infants. PsyArXiv.  
<https://doi.org/10.31234/osf.io/x9wr8>

- Rowland, C. F., Bidgood, A., Durrant, S., Peter, M., & Pine, J. M. (2018). The Language 0-5 Project. University of Liverpool. <https://doi.org/10.17605/OSF.IO/KAU5F>
- RStudio Team. (2019). *RStudio: Integrated development environment for R*. Boston, MA: RStudio, Inc. Retrieved from <http://www.rstudio.com/>
- Ryant, N., Church, K., Cieri, C., Cristia, A., Du, J., Ganapathy, S., & Liberman, M. (2019). Second DIHARD Challenge evaluation plan. *Linguistic Data Consortium, Tech. Rep.*
- Scaff, C., Stieglitz, J., Casillas, M., & Cristia, A. (n.d.). Daylong audio recordings of young children in a forager-farmer society show low levels of verbal input with minimal age-related change.
- Seidl, A., Cristia, A., Soderstrom, M., Ko, E.-S., Abel, E. A., Kellerman, A., & Schwichtenberg, A. (2018). Infant-mother acoustic-prosodic alignment and developmental risk. *Journal of Speech, Language, and Hearing Research*, 61(6), 1369–1380.
- Soderstrom, M., Bergelson, E., Warlaumont, A., Rosenberg, C., Casillas, M., Rowland, C., ... Bunce, J. (n.d.). The ACLEW Random Sampling corpus.
- Team, R. C., & others. (2013). R: A language and environment for statistical computing. Vienna, Austria.
- VanDam, M., & De Palma, P. (2018). A modular, extensible approach to massive ecologically valid behavioral data. *Behavior Research Methods*. <https://doi.org/10.3758/s13428-018-1167-8>
- VanDam, M., Warlaumont, A. S., Bergelson, E., Cristia, A., Soderstrom, M., De Palma, P., & MacWhinney, B. (2016). HomeBank: An online repository of daylong child-centered audio recordings. In *Seminars in speech and language* (Vol. 37, pp.

- 128–142). Thieme Medical Publishers.
- Warlaumont, A., Pretzer, G., Walle, E., Mendoza, S., & Lopez, L. (2016). Warlaumont HomeBank corpus.
- Weisleder, A., & Fernald, A. (2013). Talking to children matters. *Psychological Science*, *24*(11), 2143–2152. <https://doi.org/10.1177/0956797613488145>
- Xu, D., Yapanel, U., & Gray, S. (2009). Reliability of the LENATM Language Environment Analysis System in young children’s natural home environment. LENA Foundation.
- Zimmerman, F. J., Gilkerson, J., Richards, J. A., Christakis, D. A., Xu, D., Gray, S., & Yapanel, U. (2009). Teaching by listening: The importance of adult-child conversations to language development. *Pediatrics*, *124*(1), 342–349.