

¹ Peekbank: Exploring children's word recognition through an open, large-scale repository for
² developmental eye-tracking data

³ Peekbank team, Martin Zettersten¹, Claire Bergey², Naiti S. Bhatt³, Veronica Boyce⁴, Mika
⁴ Braginsky⁵, Alexandra Carstensen⁴, Benny deMayo¹, George Kachergis⁴, Molly Lewis⁶, Bria
⁵ Long⁴, Kyle MacDonald⁷, Jessica Mankewitz⁴, Stephan Meylan^{5,8}, Annissa N. Saleh⁹, Rose
⁶ M. Schneider¹⁰, Angeline Sin Mei Tsui⁴, Sarp Uner⁸, Tian Linger Xu¹¹, Daniel Yurovsky⁶, &
⁷ Michael C. Frank¹

⁸ ¹ Dept. of Psychology, Princeton University

⁹ ² Dept. of Psychology, University of Chicago

¹⁰ ³ Scripps College

¹¹ ⁴ Dept. of Psychology, Stanford University

¹² ⁵ Dept. of Brain and Cognitive Sciences, MIT

¹³ ⁶ Dept. of Psychology, Carnegie Mellon University

¹⁴ ⁷ Core Technology, McD Tech Labs

¹⁵ ⁸ Dept. of Psychology and Neuroscience, Duke University

¹⁶ ⁹ Dept. of Psychology, UT Austin

¹⁷ ¹⁰ Dept. of Psychology, UC San Diego

¹⁸ ¹¹ Dept. of Psychological and Brain Sciences, Indiana University

19

Author Note

20 Correspondence concerning this article should be addressed to Martin Zettersten.

21 E-mail: martincz@princeton.edu

22

Abstract

23 The ability to rapidly recognize words and link them to referents in context is central to
24 children's early language development. This ability, often called word recognition in the
25 developmental literature, is typically studied in the looking-while-listening paradigm, which
26 measures infants' fixation on a target object (vs. a distractor) after hearing a target label.
27 We present a large-scale, open database of infant and toddler eye-tracking data from
28 looking-while-listening tasks. The goal of this effort is to address theoretical and
29 methodological challenges in measuring vocabulary development. We first present the
30 framework for creating the database and associated tools for processing and accessing infant
31 eye-tracking datasets. Next, we show how researchers can use Peekbank to interrogate
32 theoretical and methodological questions using two illustrative examples. First, we
33 demonstrate how Peekbank can be used to investigate item-specific changes in word
34 recognition. Second, we illustrate how Peekbank can be used to create reproducible analysis
35 pipelines and to teach transparent analytic practices in infant eye-tracking research.

36 *Keywords:* word recognition; eye-tracking; vocabulary development;

37 looking-while-listening; visual world paradigm; lexical processing

38 Word count: X

39 Peekbank: Exploring children's word recognition through an open, large-scale repository for
40 developmental eye-tracking data

41 Across their first years of life, children learn words at an accelerating pace (Frank,
42 Braginsky, Yurovsky, & Marchman, 2021). While many children will only produce their first
43 word at around one year of age, most children show signs of understanding many common
44 nouns (e.g., *mommy*) and phrases (e.g., *Let's go bye-bye!*) much earlier in development
45 (Bergelson & Swingley, 2012). Although early word understanding is an enticing research
46 target, the processes involved are less directly apparent in children's behaviors and are less
47 accessible to observation than developments in speech production (Fernald, Zangl, Portillo,
48 & Marchman, 2008). To understand a spoken word, children must process the incoming
49 auditory signal and link that signal to relevant meanings – a process often referred to as
50 word recognition. A primary means of measuring word recognition in young infants are
51 eye-tracking techniques that use patterns of preferential looking to make inferences about
52 children's word processing (Fernald, Zangl, Portillo, & Marchman, 2008). The key idea of
53 these methods is that if a child preferentially looks at a target referent (rather than a
54 distractor stimulus) upon hearing a word, this indicates that the child is able to recognize
55 the word and activate its meaning during real-time language processing. Measuring early
56 word recognition offers insight into children's early word representations: children's speed of
57 response (i.e., moving their eyes; turning their heads) to the unfolding speech signal can
58 reveal children's level of comprehension (Bergelson, 2020; Fernald, Pinto, Swingley,
59 Weinberg, & McRoberts, 1998). Word recognition skills are also thought to build a
60 foundation for children's subsequent language development. Past research has found that
61 early word recognition efficiency is predictive of later linguistic and general cognitive
62 outcomes (Bleses, Makransky, Dale, Højen, & Ari, 2016; Marchman et al., 2018).

63 While word recognition is a central part of children's language development, mapping
64 the trajectory of word recognition skills has remained elusive. Studies investigating children's

word recognition are typically limited in scope to experiments in individual labs involving small samples tested on a handful of items. The limitations of single datasets makes it difficult to understand developmental changes in children's word knowledge at a broad scale. One way to overcome this challenge is to compile existing datasets into a large-scale database in order to expand the scope of research questions that can be asked about the development word recognition abilities. This strategy capitalizes on the fact that the looking-while-listening paradigm is widely used, and vast amounts of data have been collected across labs on infants' word recognition over the past 35 years (Golinkoff, Ma, Song, & Hirsh-Pasek, 2013). Such datasets have largely remained isolated from one another, but once combined, they have the potential to offer insights into the lexical development at a broad scale. Similar efforts in language development have born fruit in recent years. For example, WordBank aggregated data from the MacArthur-Bates Communicative Development Inventory, a parent-report measure of child vocabulary, to deliver new insights into cross-linguistic patterns and variability in vocabulary development (Frank, Braginsky, Yurovsky, & Marchman, 2017, 2021). In this paper, we introduce *Peekbank*, an open database of infant and toddler eye-tracking data aimed at facilitating the study of developmental changes in children's word knowledge and recognition speed.

82 The “Looking-While-Listening” Paradigm

Word recognition is traditionally studied in the “looking-while-listening” paradigm (Fernald, Zangl, Portillo, & Marchman, 2008; alternatively referred to as the intermodal preferential looking procedure, Hirsh-Pasek, Cauley, Golinkoff, & Gordon, 1987). In such studies, infants listen to a sentence prompting a specific referent (e.g., *Look at the dog!*) while viewing two images on the screen (e.g., an image of a dog – the target image – and an image of a bird – the distractor image). Infants' word recognition is measured in terms of how quickly and accurately they fixate on the correct target image after hearing its label. Past research has used this same basic method to study a wide range of questions in

language development. For example, the looking-while-listening paradigm has been used to investigate early noun knowledge, phonological representations of words, prediction during language processing, and individual differences in language development (Bergelson & Swingley, 2012; Golinkoff, Ma, Song, & Hirsh-Pasek, 2013; Lew-Williams & Fernald, 2007; Marchman et al., 2018; Swingley & Aslin, 2002).

Measuring developmental change in word recognition

While the looking-while-listening paradigm has been fruitful in advancing understanding of early word knowledge, fundamental questions remain. One central question is how to accurately capture developmental change in the speed and accuracy of word recognition. There is ample evidence demonstrating that infants get faster and more accurate in word recognition over the first few years of life (e.g., Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998). However, precisely measuring developmental increases in the speed and accuracy of word recognition remains challenging due to the difficulty of distinguishing developmental changes in word recognition skill from changes in knowledge of specific words. This problem is particularly thorny in studies with young children, since the number of items that can be tested within a single session is limited and items must be selected in an age-appropriate manner (Peter et al., 2019). Another potential challenge are that differences in the design choices and analytic decisions within single studies could obscure changes when comparing individual studies at different developmental time points. One approach to addressing these challenges is to conduct meta-analyses aggregating effects across studies while testing for heterogeneity due to researcher choices [Lewis et al. (2016); bergmann2018]. However, meta-analyses typically lack the granularity to estimate participant-level and item-level variation or to model behavior beyond coarse-grained effect size estimates. An alternative way to approach this challenge is to aggregate trial-level data from smaller studies measuring word recognition with a wide range of items and design choices into a large-scale dataset that can be analyzed using a unified modeling approach. A

117 sufficiently large dataset would allow researchers to estimate developmental change in word
118 recognition speed and accuracy while generalizing across changes related to specific words or
119 the design features of particular studies.

120 A related open theoretical question is understanding changes in children's word
121 recognition at the level of individual items. Looking-while-listening studies have been limited
122 in their ability to assess the development of specific words. One limitation is that studies
123 typically test only a small number of trials for each item, limiting the power the accurately
124 measure the development of word-specific accuracy (DeBolt, Rhemtulla, & Oakes, 2020). A
125 second limitation is that targets are often yoked with a limited set of distractors (often one
126 or two), leaving ambiguous whether accurate looking to a particular target word is largely a
127 function of children's recognition of the target word, their knowledge about the distractor,
128 which allows them to reject the distractor as a response candidate, or both. Aggregating
129 across many looking-while-listening studies has the potential to meet these challenges by
130 increasing the number of observations for specific items at different ages and by increasing
131 the variability in the distractor items co-occurring with a specific target.

132 Replicability and Reproducibility

133 A core challenge facing psychology in general, and the study of infant development in
134 particular, are threats to the replicability and reproducibility of core empirical results (Frank
135 et al., 2017; Nosek et al., 2021). In infant research, many studies are not adequately powered
136 to detect the main effects of interest (Bergmann et al., 2018). This is often compounded by
137 low reliability in infant measures, often due to limits on the number of trials that can be
138 collected from an individual infant in an experimental session (Byers-Heinlein, Bergmann, &
139 Savalei, 2021). One hurdle to improving the power in infant research is that it can often be
140 difficult to develop a priori estimates of effect sizes, and how specific design decisions (e.g.,
141 the number of test trials) will impact power and reliability. Large-scale databases of infant
142 behavior can aid researchers' in their decision-making by providing rich datasets that can

143 help constrain expectations about possible effect sizes and can be used to make data-driven
144 design decisions. For example, if a researcher is interested in understanding how the number
145 of test trials could impact the power and reliability of their looking-while-listening design, a
146 large-scale database would allow them to simulate possible outcomes across a range of test
147 trials, based on past eye-tracking data with infants.

148 In addition to threats to replicability, the field of infant development also faces
149 concerns about analytic reproducibility - the ability for researchers to arrive at the same
150 analytic conclusion reported in the original research article, given the same dataset. A recent
151 estimate based on studies published in a prominent cognitive science journal suggests that
152 analyses can remain difficult to reproduce, even when data is made available to other
153 research teams (Hardwicke et al., 2018). Aggregating data in centralized databases can aid
154 in improving reproducibility in several ways. First, building a large-scale database requires
155 defining a standardized data specification. Recent examples include the brain imaging data
156 structure (BIDS), an effort to specify a unified data format for neuroimaging experiments
157 (Gorgolewski et al., 2016). Defining a data standard - in this case, for infant eye-tracking
158 experiments - supports reproducibility by setting data curation standards that guarantee
159 that critical information will be available in openly shared data and that make it easier for
160 different research teams to understand the data structure. Second, open databases make it
161 easy for researchers to generate open and reproducible analytic pipelines, both for individual
162 studies and for analyses aggregating across datasets. Creating open analytic pipelines across
163 many datasets also serves a pedagogical purpose, providing teaching examples illustrating
164 how to implement analytic techniques used in influential studies and how to conduct
165 reproducible analyses with infant eye-tracking data.

166 **Peekbank: An open database of developmental eye-tracking studies.**

167 What all of these open challenges share is that they are difficult to address at the scale
168 of a single research lab or in a single study. To address this challenge, we developed

169 *Peekbank* a flexible and reproducible interface to an open database of developmental
170 eye-tracking studies. The Peekbank project (a) collects a large set of eye-tracking datasets
171 on children’s word recognition, (b) introduces a data format and processing tools for
172 standardizing eye-tracking data across heterogeneous data sources, and (c) provides an
173 interface for accessing and analyzing the database. In the current paper, we introduce the
174 key components of the project and give an overview of the existing database. We then
175 provide worked examples of how researchers can use Peekbank (1) to inform methodological
176 decision-making, (2) to teach through reproducible examples, and (3) ask novel research
177 questions about the development of children’s word recognition.

178 **Design and Technical Approach**

179 **Database Framework**

180 One of the main challenges in compiling a large-scale eye-tracking database is the lack
181 of a shared data format: both labs and individual experiments can record their results in a
182 wide range of formats. For example, different experiments encode trial-level and subject-level
183 information in many different ways. Therefore, we have developed a common tabular format
184 to support analyses of all studies simultaneously.

185 As illustrated in Figure 1, the Peekbank framework consists of four main components:
186 (1) a set of tools to *convert* eye-tracking datasets into a unified format, (2) a relational
187 database populated with data in this unified format, (3) a set of tools to *retrieve* data from
188 this database, and (4) a web app (using the Shiny framework) for visualizing the data. These
189 components are supported by three packages. The `peekds` package (for the R language; R
190 Core Team (2020)) helps researchers convert existing datasets to use the standardized format
191 of the database. The `peekbank` module (Python) creates a database with the relational
192 schema and populates it with the standardized datasets produced by `peekds`. The database
193 is served through MySQL, an industry standard relational database server, which may be
194 accessed by a variety of programming languages, and can be hosted on one machine and

accessed by many others over the Internet. As is common in relational databases, records of similar types (e.g., participants, trials, experiments, coded looks at each timepoint) are grouped into tables, and records of various types are linked through numeric identifiers. The `peekbankr` package (R) provides an application programming interface, or API, that offers high-level abstractions for accessing the tabular data stored in Peekbank. Most users will access data through this final package, in which case the details of data formatting, processing, and the specifics of connecting to the database are abstracted away from the user.

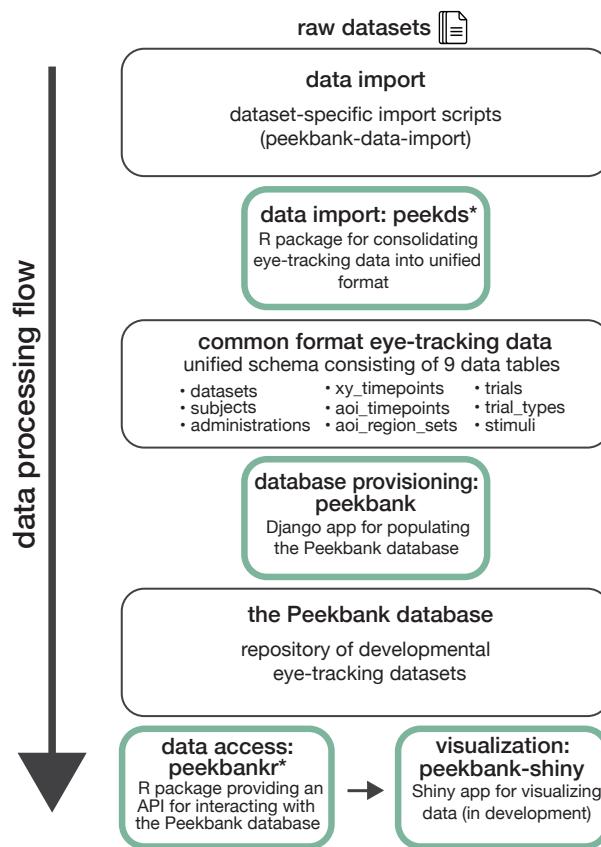


Figure 1. Overview of the Peekbank data ecosystem. Peekbank tools are highlighted in green.
* indicates R packages introduced in this work.

202 Database Schema

The Peekbank database contains two major types of data: (1) metadata regarding experiments, participants, and trials, and (2) time course looking data, detailing where on the screen a child is looking at a given point in time (Fig. 2).

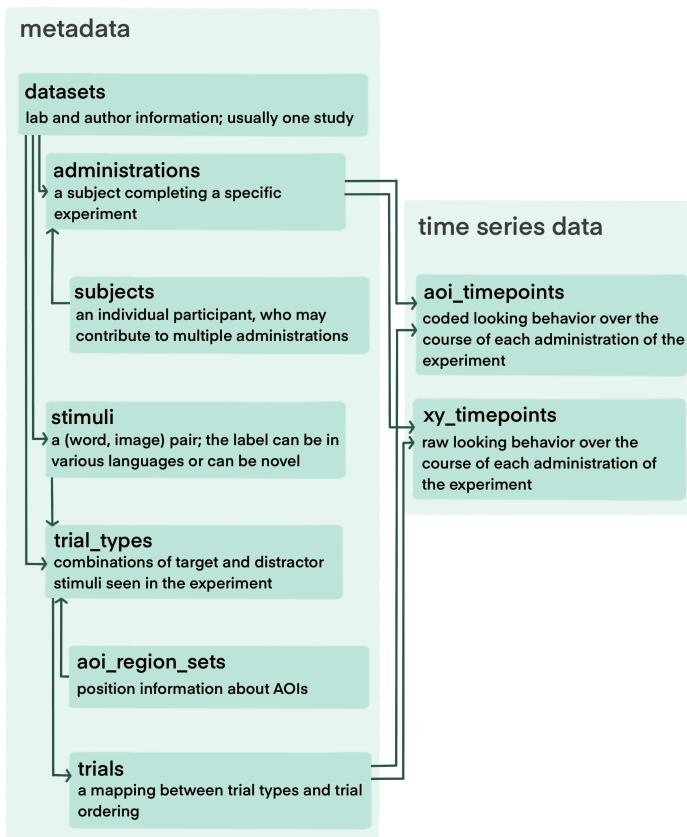


Figure 2. The Peekbank schema. Each square represents a table in the relational database.

206 **Metadata.** Metadata can be separated into four parts: (1) participant-level
 207 information (e.g., demographics) (2) experiment-level information (e.g., the type of eye
 208 tracker used to collect the data) (3) session information (e.g. a participant's age for a specific
 209 experimental session) and (4) trial information (e.g., what images or videos were presented
 210 onscreen, and paired with which audio).

211 **Participant Information.** Invariant information about individuals who
 212 participate in one or more studies (e.g, a subject's first language) is recorded in the
 213 **subjects** table, while the **administrations** table contains information about a subject's
 214 participation in a single session of a study (see Session Information, below). This division
 215 allows Peekbank to gracefully handle longitudinal designs: a single subject can be associated
 216 with many administrations.

217 Subject-level data includes all participants who have experiment data. In general, we
 218 include as many participants as possible in the database and leave it to end-users to apply
 219 the appropriate exclusion criteria for their analysis.

220 ***Experiment Information.*** The `datasets` table includes information about the
 221 lab conducting the study and the relevant publications to cite regarding the data. In most
 222 cases, a dataset corresponds to a single study.

223 Information about the experimental design is split across the `trial_types` and
 224 `stimuli` tables. The `trial_types` table encodes information about each trial *in the design*
 225 *of the experiment*,¹ including the target stimulus and location (left vs. right), the distractor
 226 stimulus and location, and the point of disambiguation for that trial. If a dataset used
 227 automatic eye-tracking rather than manual coding, each trial type is additionally linked to a
 228 set of area of interest (x, y) coordinates, encoded in the `aoi_region_sets` table. The
 229 `trial_types` table links trial types to the `aoi_region_sets` table and the `trials` table.
 230 Each trial_type record links to two records in the `stimuli` table, identified by the
 231 `distractor_id` and the `target_id` fields.

232 Each record in the `stimuli` table is a (word, image) pair. In most experiments, there is
 233 a one-to-one mapping between images and labels (e.g., each time an image of a dog appears
 234 it is referred to as *dog*). For studies in which there are multiple potential labels per image
 235 (e.g., *dog* and *chien* are both used to refer to an image of a dog), images can have multiple
 236 rows in the `stimuli` table with unique labels as well as a row with no label to be used when
 237 the image appears solely as a distractor (and thus its label is ambiguous). This structure is
 238 useful for studies on synonymy or using multiple languages. For studies in which the same
 239 label refers to multiple images (e.g., the word *dog* refers to an image of a dalmatian and a
 240 poodle), the same label can have multiple rows in the `stimuli` table with unique images.

¹ We note that the term *trial* is often overloaded, to refer to a particular combination of stimuli seen by many participants, vs. a participant seeing that particular combination at a particular point in the experiment. We track the latter in the ‘trials’ table.

241 **Session Information.** The `administrations` table includes information about

242 the participant or experiment that may change between sessions of the same study, even for

243 the same participant. This includes the age of the participant, the coding method

244 (eye-tracking vs. hand-coding), and the properties of the monitor that was used.

245 **Trial Information.** The `trials` table includes information about a specific

246 participant completing a specific instance of a trial type. This table links each record in the

247 raw data (described below) to the trial type and specifies the order of the trials seen by a

248 specific participant.

249 **Time course data.** Raw looking data is a series of looks to areas of interest (AOIs),

250 such as looks to the left or right of the screen, or to (x, y) coordinates on the experiment

251 screen, linked to points in time. For data generated by eye-trackers, we typically have (x, y)

252 coordinates at each time point, which will be encoded in the `xy_timepoints` table. These

253 looks will also be recoded into AOIs according to the AOI coordinates in the

254 `aoi_region_sets` table using the `add_aois()` function in `peekds`, which will be encoded in

255 the `aoi_timepoints` table. For hand-coded data, we typically have a series of AOIs (i.e.,

256 looks to the left vs. right of the screen), but lack information about exact gaze positions

257 on-screen; the AOIs will be recoded into the categories in the Peekbank schema (target,

258 distractor, other, and missing) and encoded in the `aoi_timepoints` table, and these

259 datasets will not have an `xy_timepoints` table.

260 Typically, timepoints in the `xy_timepoints` table and `aoi_timepoints` table need to

261 be regularized to center each trial's time around the point of disambiguation—such that 0 is

262 the time of target word onset in the trial (i.e., the beginning of *dog* in *Can you find the*

263 *dog?*). While information preceding the onset of the target label in some datasets, such as

264 coarticulation cues (Mahr, McMillan, Saffran, Ellis Weismer, & Edwards, 2015) and specific

265 adjectives (Fernald, Marchman, & Weisleder, 2013), can in principle disambiguate the target

266 referent, we re-centered timing information to the onset of the target label to facilitate

267 comparison of target label processing across all datasets. If time values run throughout the

268 experiment rather than resetting to zero at the beginning of each trial, `rezero_times()` is
269 used to reset the time at each trial. After this, each trial's times are centered around the
270 point of disambiguation using `normalize_times()`. When these steps are complete, the
271 time course is ready for resampling.

272 To facilitate time course analysis and visualization across datasets, time course data
273 must be resampled to a uniform sampling rate (i.e., such that every trial in every dataset has
274 observations at the same time points). To do this, we use the `resample_times()` function.
275 During the resampling process, we interpolate using constant interpolation, selecting for each
276 interpolated timepoint the looking location for the nearest observed time point in the
277 original data for both `aoi_timepoints` and `xy_timepoints` data. In the case of ties, the
278 look location observed at the earlier timepoint in the original data is chosen for the
279 resampled timepoint. Currently, all data is resampled to 40 Hz (observations every 25 ms) by
280 default, which represents a compromise between retaining fine-grained timing information
281 from datasets with dense sampling rates (maximum sampling rate among current datasets:
282 500 Hz) while minimizing the possibility of introducing artifacts via resampling for datasets
283 with lower sampling rates (minimum sampling rate for current datasets: 30 Hz). Compared
284 to linear interpolation (see e.g. Wass et al., 2014), constant interpolation has the advantage
285 that it is more conservative, in the sense that it does not introduce new look locations
286 beyond those measured in the original data.

287 Processing, Validation and Ingestion

288 The `peekds` package offers functions to extract the above data. Once this data has
289 been extracted in a tabular form, the package also offers a function to check whether all
290 tables have the required fields and data types expected by the database. In an effort to
291 double check the data quality and to make sure that no errors are made in the importing
292 script, as part of the import procedure we create a time course plot based on our processed
293 tables to replicate the results in the paper that first presented each dataset. Once this plot

294 has been created and checked for consistency and all tables pass our validation functions, the
 295 processed dataset is ready for reprocessing into the database using the `peekbank` library.
 296 This library applies additional data checks, and adds the data to the MySQL database using
 297 the Django web framework.

298 Currently, the import process is carried out by the Peekbank team using data offered
 299 by other research teams. In the future, we hope to allow research teams to carry out their
 300 own import processes with checks from the Peekbank team before reprocessing. To this end,
 301 import script templates are available for both hand-coded datasets and automatic
 302 eye-tracking datasets for research teams to adapt to their data.

303 **Current Data Sources**

Table 1
Overview of the datasets in the current database.

Dataset name	Citation	N	Mean age (mos.)	Age range (mos.)	Method	Language
attword	Yurovsky & Frank, 2017	288	25.5	13–59	eye-tracking	English
canine	unpublished	36	23.8	21–27	manual coding	English
coartic	Mahr et al., 2015	29	20.8	18–24	eye-tracking	English
cowpig	Perry et al., 2017	45	20.5	19–22	manual coding	English
fmw	Fernald et al., 2013	80	20.0	17–26	manual coding	English
ft_pt	Adams et al., 2018	69	17.1	13–20	manual coding	English
input_uptake	Hurtado et al., 2008	76	21.0	17–27	manual coding	Spanish
lsc	Ronfard et al., 2021	40	20.0	18–24	manual coding	English
mispron	Swingley & Aslin, 2002	50	15.1	14–16	manual coding	English
mix	Byers-Heinlein et al., 2017	48	20.1	19–21	eye-tracking	English, French
reflook_socword	Yurovsky et al., 2013	435	33.6	12–70	eye-tracking	English
reflook_v4	unpublished	45	34.2	11–60	eye-tracking	English
remix	Potter et al., 2019	44	22.6	18–29	manual coding	Spanish, English
salientme	Pomper & Saffran, 2019	44	40.1	38–43	manual coding	English
stl	Weisleder & Fernald, 2013	29	21.6	18–27	manual coding	Spanish
switchingCues	Pomper & Saffran, 2016	60	44.3	41–47	manual coding	English
tablet	Frank et al., 2016	69	35.5	12–60	eye-tracking	English
tseltal	Casillas et al., 2017	23	31.3	9–48	manual coding	Tseltal
xsectional	Hurtado et al., 2007	49	23.8	15–37	manual coding	Spanish
yoursmy	Garrison et al., 2020	35	14.5	12–18	eye-tracking	English

304 The database currently includes 20 looking-while-listening datasets comprising $N=1594$
 305 total participants (Table 1). The current data represents a convenience sample of datasets
 306 that were (a) datasets collected by or available to Peekbank team members, (b) made
 307 available to Peekbank after informal inquiry or (c) datasets that were openly available. Most

308 datasets (14 out of 20 total) consist of data from monolingual native English speakers. They
309 span a wide age spectrum with participants ranging from 9 to 70 months of age, and are
310 balanced in terms of gender (47% female). The datasets vary across a number of
311 design-related dimensions, and include studies using manually coded video recordings and
312 automated eye-tracking methods (e.g., Tobii, EyeLink) to measure gaze behavior. All studies
313 tested familiar items, but the database also includes 5 datasets that tested novel
314 pseudo-words in addition to familiar words.

315 **Versioning + Expanding the database**

316 The content of Peekbank will change as we add additional datasets and revise previous
317 ones. To facilitate reproducibility of analyses, we use a versioning system where successive
318 releases are assigned a name reflecting the year and version, e.g., 2021.1. By default, users
319 will interact with the most recent version of the database available, though `peekbankr` API
320 allows researchers to run analyses against any previous version of the database. For users
321 with intensive use-cases, each version of the database may be downloaded as a compressed
322 .sql file and installed on a local MySQL server.

323 **Interfacing with peekbank**

324 **Peekbankr**

325 The `peekbankr` API offers a way for users to access data from the database and
326 flexibly analyze it in R. Users can download tables from the database, as specified in the
327 Schema section above, and merge them using their linked IDs to examine time course data
328 and metadata jointly. In the sections below, we work through some examples to outline the
329 possibilities for analyzing data downloaded using `peekbankr`.

330 Functions:

- 331 • `connect_to_peekbank()` opens a connection with the Peekbank database to allow

332 tables to be downloaded with the following functions

- 333 • `get_datasets()` gives each dataset name and its citation information
- 334 • `get_subjects()` gives information about persistent subject identifiers (e.g., native
languages, sex)
- 335 • `get_administrations()` gives information about specific experimental
administrations (e.g., subject age, monitor size, gaze coding method)
- 336 • `get_stimuli()` gives information about word–image pairings that appeared in
experiments
- 337 • `get_trial_types()` gives information about pairings of stimuli that appeared in the
experiment (e.g., point of disambiguation, target and distractor stimuli, condition,
language)
- 338 • `get_trials()` gives the trial orderings for each administration, linking trial types to
the trial IDs used in time course data
- 339 • `get_aoi_region_sets()` gives coordinate regions for each area of interest (AOI)
linked to trial type IDs
- 340 • `get_xy_timepoints()` gives time course data for each subject’s looking behavior in
each trial, as (x, y) coordinates on the experiment monitor
- 341 • `get_aoi_timepoints()` gives time course data for each subject’s looking behavior in
each trial, coded into areas of interest

351 Shiny App

352 One goal of the Peekbank project is to allow a wide range of users to easily explore and

353 learn from the database. We therefore have created an interactive web application –

354 `peekbank-shiny` – that allows users to quickly and easily create informative visualizations

355 of individual datasets and aggregated data. `peekbank-shiny` is built using Shiny, a software

356 package for creating web apps for data exploration with R, as well as the `peekbankr` package.

357 The Shiny app allows users to create commonly used visualizations of looking-while-listening

358 data, based on data from the Peekbank database. Specifically, users can visualize
359 1. the time course of looking data in a profile plot depicting infant target looking across
360 trial time
361 2. overall accuracy (proportion target looking) within a specified analysis window
362 3. reaction times (speed of fixating the target image) in response to a target label
363 4. an onset-contingent plot, which shows the time course of participant looking as a
364 function of their look location at the onset of the target label

365 Users are given various customization options for each of these visualizations, e.g.,
366 choosing which datasets to include in the plots, controlling the age range of participants,
367 splitting the visualizations by age bins, and controlling the analysis window for time course
368 analyses. Plots are then updated in real time to reflect users' customization choices, and
369 users are given options to share the visualizations they created. An screenshot of the app is
370 shown in Figure ???. The Shiny app thus allows users to quickly inspect basic properties of
371 Peekbanks datasets and create reproducible visualizations without incurring any of the
372 technical overhead required to access the database through R.

373 OSF site

374 In addition to the Peekbank database proper, all data is openly available on the
375 Peekbank OSF webpage (<https://osf.io/pr6wu/>). The OSF site also includes the original raw
376 data (both time series data and metadata, such as trial lists and participant logs) that was
377 obtained for each study and subsequently processed into the standardized Peekbank format.
378 Users who are interested in inspecting or reproducing the processing pipeline for a given
379 dataset can use the respective import script (openly available on GitHub,
380 <https://github.com/langcog/peekbank-data-import>) to download and process the raw data
381 from OSF into its final standardized format. Where available, the OSF page also includes
382 additional information about the stimuli used in each dataset, including in some instances

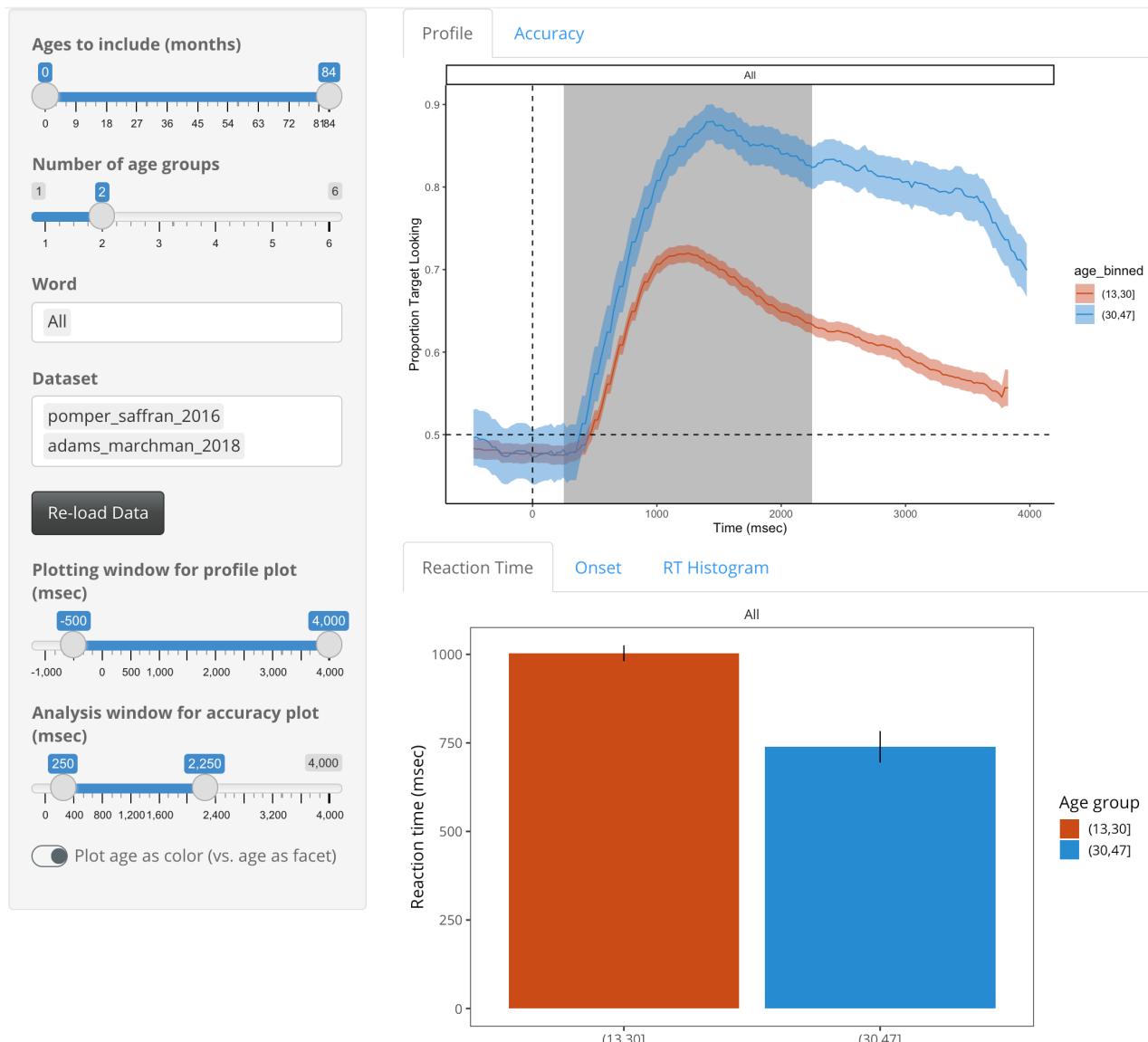


Figure 3. Screenshot of the Peekbank visualization tool, which shows a variety of standard analysis plots as a function of user-selected datasets, words, age ranges, and analysis windows. Shown here are mean reaction time and proportion target looking over time by age group for two selected datasets.

383 the original stimulus sets (e.g., image and audio files).

384 **Peekbank: General Descriptives**

385 [Accuracy, Reaction Times, Item variability?]

386 **Overall Word Recognition Accuracy**

Dataset Name	Unique Items	Prop. Target	95% CI
attword	6	0.63	[0.62, 0.65]
canine	16	0.65	[0.61, 0.68]
coartic	10	0.71	[0.68, 0.74]
cowpig	12	0.61	[0.58, 0.63]
fmw	12	0.65	[0.63, 0.67]
ft_pt	8	0.65	[0.63, 0.67]
input_uptake	12	0.61	[0.59, 0.63]
lsc	8	0.69	[0.65, 0.73]
mispron	22	0.57	[0.55, 0.59]
mix	6	0.55	[0.52, 0.58]
reflook_socword	6	0.61	[0.6, 0.63]
reflook_v4	10	0.61	[0.57, 0.65]
remix	8	0.63	[0.58, 0.67]
salientme	16	0.74	[0.72, 0.75]
stl	12	0.63	[0.6, 0.66]
switchingCues	40	0.77	[0.75, 0.8]
tablet	24	0.64	[0.6, 0.68]
tseltal	30	0.59	[0.54, 0.63]
xsectional	8	0.59	[0.55, 0.63]
yoursmy	87	0.60	[0.56, 0.64]

Table 2

Average proportion target looking in each dataset.

387 In general, participants demonstrated robust, above-chance word recognition in each
 388 dataset (chance=0.5). Table 2 shows the average proportion of target looking within a
 389 standard critical window of 367-2000ms after the onset of the label for each dataset
 390 (Swingley & Aslin, 2002). Proportion target looking was generally higher for familiar words
 391 ($M = 0.66$, 95% CI = [0.65, 0.67], $n = 1543$) than for novel words learned during the
 392 experiment ($M = 0.59$, 95% CI = [0.58, 0.61], $n = 822$).

393 **Item-level variability**

394 Figure 4 gives an overview of the variability in accuracy for individual words in each
 395 dataset. The number of unique target labels and their associated accuracy vary widely

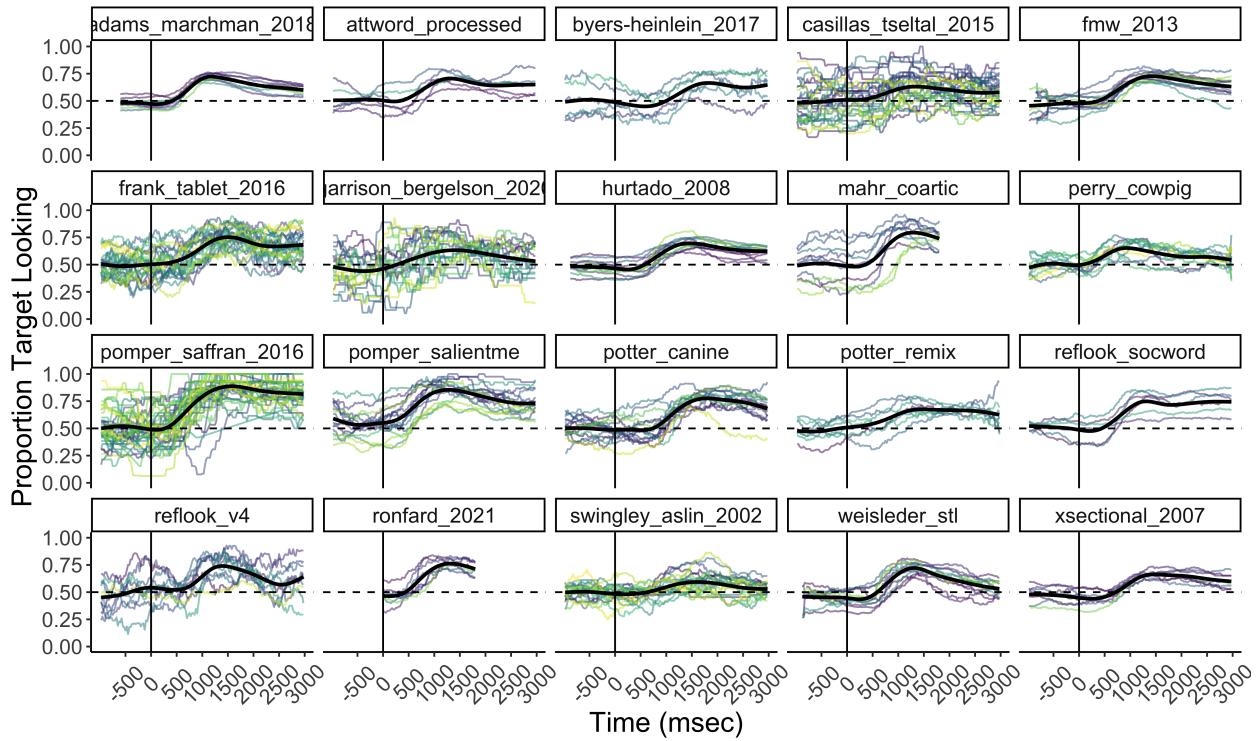


Figure 4. Item-level variability in proportion target looking within each dataset (chance=0.5). Time is centered on the onset of the target label (vertical line). Colored lines represent specific target labels. Black lines represent smoothed average fits based on a general additive model using cubic splines.

396 across datasets.

397

Peekbank in Action

398 We provide two potential use-cases for Peekbank data. In each case, we provide sample
 399 code so as to model how easy it is to do simple analyses using data from the database. Our
 400 first example shows how we can replicate the analysis for a classic study. This type of
 401 computational reproducibility can be a very useful exercise for teaching students about best
 402 practices for data analysis (e.g., Hardwicke et al., 2018) and also provides an easy way to
 403 explore looking-while-listening time course data in a standardized format. Our second
 404 example shows an in-depth exploration of developmental changes in the recognition of
 405 particular words. Besides its theoretical interest (which we will explore more fully in
 406 subsequent work), this type of analysis could in principle be used for optimizing the stimuli

407 for new experiments, especially as the Peekbank dataset grows and gains coverage over a
 408 greater number of items.

409 **Computational reproducibility example: Swingley and Aslin (2002)**

410 Swingley and Aslin (2002) investigated the specificity of 14-16 month-olds' word
 411 representations using the looking-while-listening paradigm, asking whether recognition would
 412 be slower and less accurate for mispronunciations, e.g. *oppel* (close mispronunciation) or *opel*
 413 (distant mispronunciation) instead of *apple* (correct pronunciation). In this short vignette,
 414 we show how easily the data in Peekbank can be used to visualize this result.

```
library(peekbankr)
aoi_timepoints <- get_aoi_timepoints(dataset_name = "swingley_aslin_2002")
administrations <- get_administrations(dataset_name = "swingley_aslin_2002")
trial_types <- get_trial_types(dataset_name = "swingley_aslin_2002")
trials <- get_trials(dataset_name = "swingley_aslin_2002")
```

415 We begin by retrieving the relevant tables from the database, `aoi_timepoints`,
 416 `administrations`, `trial_types`, and `trials`. As discussed above, each of these can be
 417 downloaded using a simple API call through `peekbankr`, which returns dataframes that
 418 include ID fields. These ID fields allow for easy joining of the data into a single dataframe
 419 containing all the information necessary for the analysis.

```
swingley_data <- aoi_timepoints %>%
  left_join(administrations) %>%
  left_join(trials) %>%
  left_join(trial_types) %>%
  filter(condition != "filler") %>%
  mutate(condition = if_else(condition == "cp", "Correct", "Mispronounced"))
```

420 As the code above shows, once the data are joined, condition information for each
 421 timepoint is present and so we can easily filter out filler trials and set up the conditions for
 422 further analysis. For simplicity, here we combine both mispronunciation conditions since the
 423 close vs. distant mispronunciation manipulation showed no effect in the original paper.

```

accuracies <- swingley_data %>%
  group_by(condition, t_norm, administration_id) %>%
  summarize(correct = sum(aoi == "target") /
    sum(aoi %in% c("target", "distractor"))) %>%
  group_by(condition, t_norm) %>%
  summarize(mean_correct = mean(correct),
    ci = 1.96 * sd(correct) / sqrt(n()))

```

424 The final step in our analysis is to create a summary dataframe using `dplyr`

425 commands. We first group the data by timestep, participant, and condition and compute the
 426 proportion looking at the correct image. We then summarize again, averaging across
 427 participants, computing both means and 95% confidence intervals (via the approximation of
 428 1.96 times the standard error of the mean). The resulting dataframe can be used for
 429 visualization of the time course of looking.

430 Figure 5 shows the average time course of looking for the two conditions, as produced

431 by the code above. Looks after the correctly pronounced noun appeared both faster
 432 (deviating from chance earlier) and more accurate (showing a higher asymptote). Overall,
 433 this example demonstrates the ability to produce this visualization in just a few lines of code.

434 Item analyses

435 A second use case for Peekbank is to examine item-level variation in word recognition.

436 Individual datasets rarely have enough statistical power to show reliable developmental
 437 differences within items. To illustrate the power of aggregating data across multiple datasets,
 438 we select the four words with the most data available across studies and ages (apple, book,
 439 dog, and frog) and show average recognition trajectories.

440 Our first step is to collect and join the data from the relevant tables including

441 timepoint data, trial and stimulus data, and administration data (for participant ages). We
 442 join these into a single dataframe for easy manipulation; this dataframe is a common
 443 starting point for analyses of item-level data.

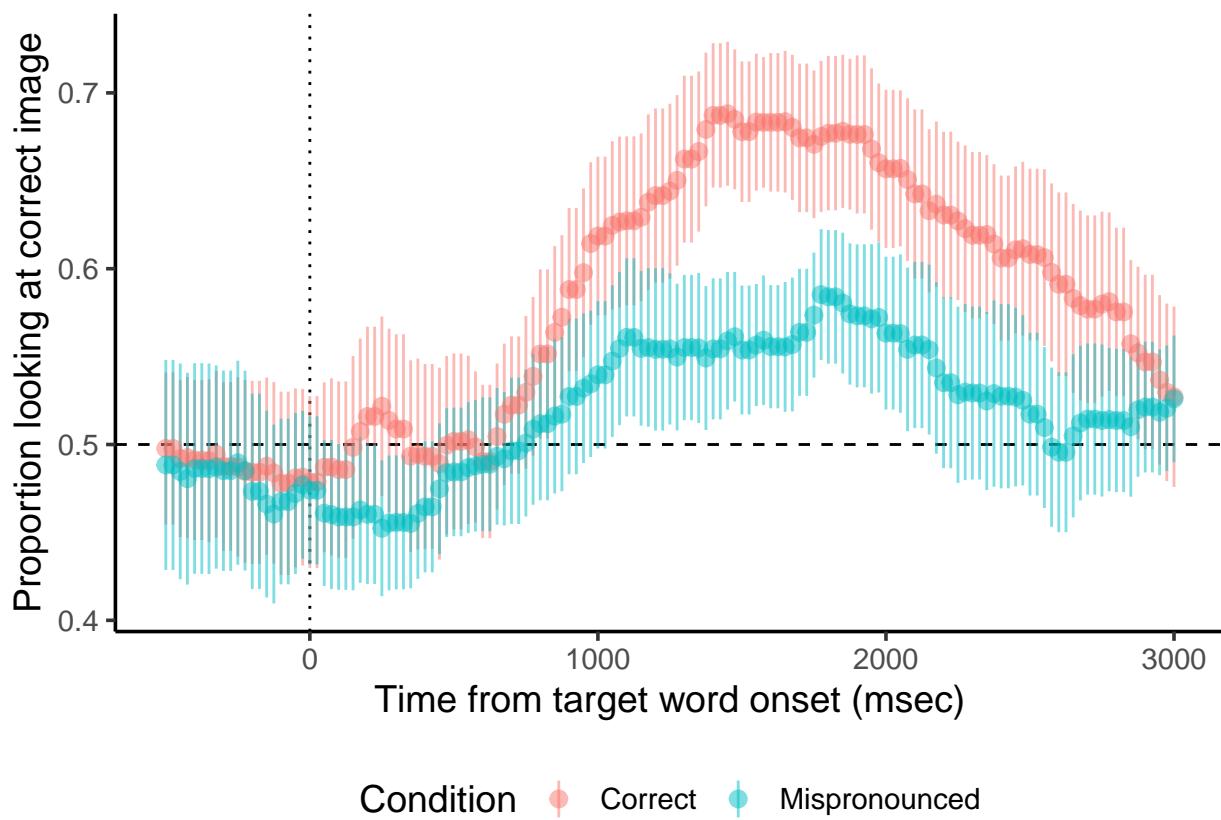


Figure 5. Proportion looking at the correct referent by time from the point of disambiguation (the onset of the target noun) in Swoley & Aslin (2002). Colors show the two pronunciation conditions; points give means and ranges show 95% confidence intervals. The dotted line shows the point of disambiguation and the dashed line shows chance performance.

```

all_aoi_timepoints <- get_aoi_timepoints()
all_stimuli <- get_stimuli()
all_administrations <- get_administrations()
all_trial_types <- get_trial_types()
all_trials <- get_trials()

aoi_data_joined <- all_aoi_timepoints %>%
  right_join(all_administrations) %>%
  right_join(all_trials) %>%
  right_join(all_trial_types) %>%

```

```

  mutate(stimulus_id = target_id) %>%
  right_join(all_stimuli) %>%
  select(administration_id, english_stimulus_label, age, t_norm, aoi)

```

444 Next we select a set of four target words (chosen based on having more than XXX
 445 children contributing data for each across several one-year age groups). We create age
 446 groups, aggregate, and compute timepoint-by-timepoint confidence intervals using the z
 447 approximation.

```

target_words <- c("book", "dog", "frog", "apple")

target_word_data <- aoi_data_joined %>%
  filter(english_stimulus_label %in% target_words) %>%
  mutate(age_group = cut(age, breaks = seq(12, 48, 12))) %>%
  filter(!is.na(age_group)) %>%
  group_by(t_norm, administration_id, age_group, english_stimulus_label) %>%
  summarise(correct = mean(aoi == "target") /
    mean(aoi %in% c("target", "distractor"), na.rm=TRUE)) %>%
  group_by(t_norm, age_group, english_stimulus_label) %>%
  summarise(ci = 1.96 * sd(correct, na.rm=TRUE) / sqrt(length(correct)),
    correct = mean(correct, na.rm=TRUE),
    n = n())

```

448 Finally, we plot the data as time courses split by age. Our plotting code is shown
 449 below (with styling commands again removed for clarity). Figure 6 shows the resulting plot,
 450 with time courses for each of three (rather coarse) age bins. Although some baseline effects
 451 are visible across items, we still see clear and consistent increases in looking to the target,
 452 with the increase appearing earlier and in many cases asymptoting at a higher level for older

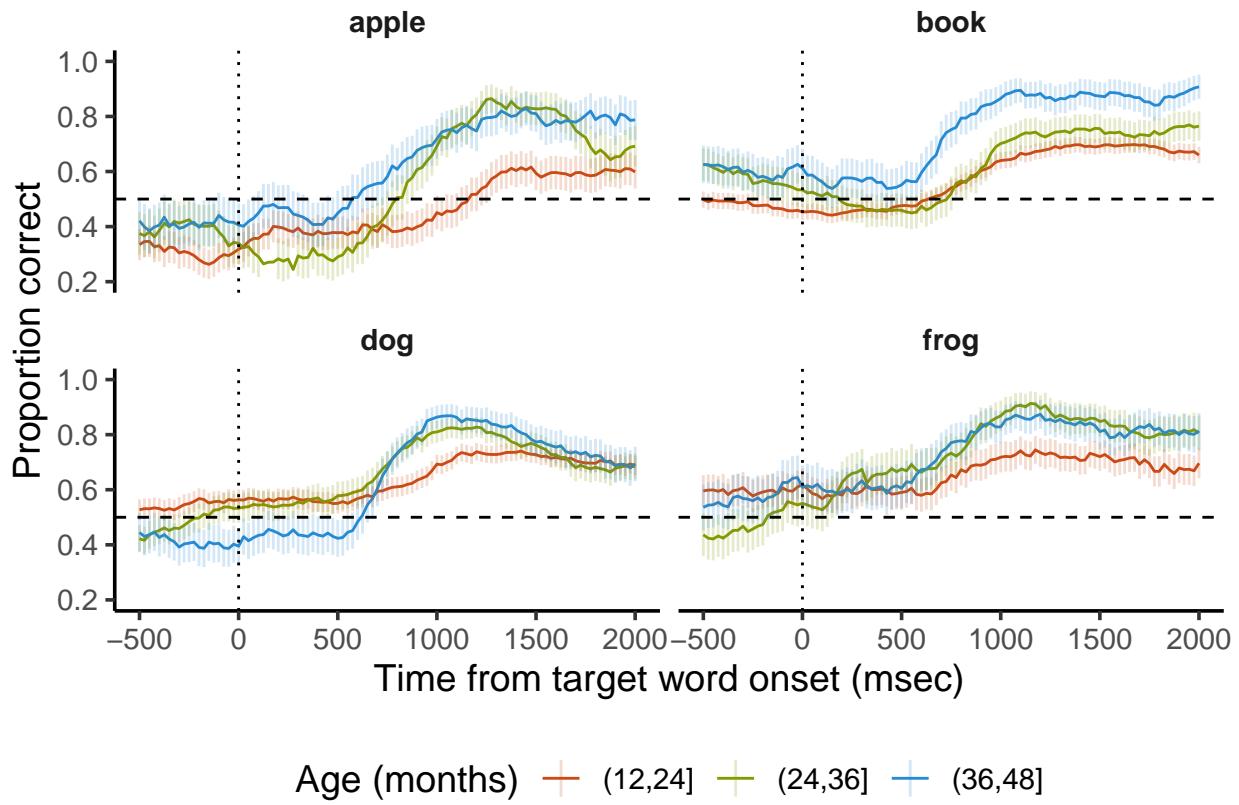


Figure 6. Time course plot for four well-represented target items in the Peekbank dataset, split by three age groups. Each line represents children's average looking to the target image after the onset of the target label (dashed vertical line). Error bars represent 95% CIs.

453 children. On the other hand, this simple averaging approach ignores study-to-study variation
 454 (perhaps responsible for the baseline effects we see in the *apple* and *frog* items especially). In
 455 future work, we hope to introduce model-based analytic methods that use mixed effects
 456 regression to factor out study-level and individual-level variance in order to recover
 457 developmental effects more appropriately (see e.g. Zettersten et al. (2021) for a prototype of
 458 such an analysis).

```
ggplot(target_word_data,
       aes(x = t_norm, y = correct, col = age_group)) +
  geom_line() +
  geom_linerange(aes(ymin = correct - ci, ymax = correct + ci),
                 alpha = .2) +
```

```
facet_wrap(~english_stimulus_label)
```

459

Discussion

460 Theoretical progress in understanding child development requires rich datasets, but
 461 collecting child data is expensive, difficult, and time-intensive. Recent years have seen a
 462 growing effort to build open source tools and pool research efforts to meet the challenge of
 463 building a cumulative developmental science (Bergmann et al., 2018; Frank, Braginsky,
 464 Yurovsky, & Marchman, 2017; Sanchez et al., 2019; The ManyBabies Consortium, 2020).
 465 The Peekbank project expands on these efforts by building an infrastructure for aggregating
 466 eye-tracking data across studies, with a specific focus on the looking-while-listening
 467 paradigm. This paper presents an overview of the structure of the database, as well as how
 468 users can access the database and some initial demonstrations of how it can be used both to
 469 facilitate reproducibility, for teaching and for exploring theoretical questions beyond on the
 470 scope of an individual study.

471 There are a number of limitations surrounding the current scope of the database. A
 472 priority in future work will be to expand the size of the database. With 20 datasets currently
 473 available in the database, idiosyncrasies of particular designs and condition manipulations
 474 still have substantial influence on modeling results. Expanding the set of distinct datasets
 475 will allow us to increase the number of observations per item across datasets, leading to more
 476 robust generalizations across item-level variability. The current database is also limited by
 477 the relatively homogeneous background of its participants, both with respect to language
 478 (almost entirely monolingual native English speakers) and cultural background (Henrich,
 479 Heine, & Norenzayan, 2010; Muthukrishna et al., 2020). Increasing the diversity of
 480 participant backgrounds and languages will expand the scope of the generalizations we can
 481 form about child word recognition.

482 Finally, while the current database is focused on studies of word recognition, the tools

483 and infrastructure developed in the project can in principle be used to accommodate any
484 eye-tracking paradigm, opening up new avenues for insights into cognitive development. Gaze
485 behavior has been at the core of many of the key advances in our understanding of infant
486 cognition. Aggregating large datasets of infant looking behavior in a single, openly-accessible
487 format promises to bring a fuller picture of infant cognitive development into view.

488 **Acknowledgements**

489 We would like to thank the labs and researchers that have made their data publicly
490 available in the database.

491

References

- 492 Bergelson, E. (2020). The comprehension boost in early word learning: Older infants
493 are better learners. *Child Development Perspectives*, 14(3), 142–149.
- 494 Bergelson, E., & Swingley, D. (2012). At 6-9 months, human infants know the
495 meanings of many common nouns. *PNAS*, 109(9), 3253–3258.
- 496 Bergmann, C., Tsuji, S., Piccinini, P. E., Lewis, M. L., Braginsky, M., Frank, M. C.,
497 & Cristia, A. (2018). Promoting replicability in developmental research through
498 meta-analyses: Insights from language acquisition research. *Child Development*,
499 89(6), 1996–2009.
- 500 Bleses, D., Makransky, G., Dale, P. S., Højen, A., & Ari, B. A. (2016). Early
501 productive vocabulary predicts academic achievement 10 years later. *Applied
502 Psycholinguistics*, 37(6), 1461–1476.
- 503 Byers-Heinlein, K., Bergmann, C., & Savalei, V. (2021). Six solutions for more reliable
504 infant research. *PsyArXiv*. <https://doi.org/https://doi.org/10.31234/osf.io/ksfvq>
- 505 DeBolt, M. C., Rhemtulla, M., & Oakes, L. M. (2020). Robust data and power in
506 infant research: A case study of the effect of number of infants and number of
507 trials in visual preference procedures. *Infancy*, 25(4), 393–419.
508 <https://doi.org/10.1111/infa.12337>
- 509 Fernald, A., Marchman, V. A., & Weisleder, A. (2013). SES differences in language
510 processing skill and vocabulary are evident at 18 months. *Developmental Science*,
511 16(2), 234–248. <https://doi.org/10.1111/desc.12019>
- 512 Fernald, A., Pinto, J. P., Swingley, D., Weinberg, A., & McRoberts, G. W. (1998).
513 Rapid gains in speed of verbal processing by infants in the 2nd year. *Psychological*

514 *Science*, 9(3), 228–231.

515 Fernald, A., Zangl, R., Portillo, A. L., & Marchman, V. A. (2008). Looking while
516 listening: Using eye movements to monitor spoken language comprehension by
517 infants and young children. In I. A. Sekerina, E. M. Fernandez, & H. Clahsen
518 (Eds.), *Developmental psycholinguistics: On-line methods in children's language*
519 *processing* (pp. 97–135). Amsterdam: John Benjamins.

520 Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., ...
521 Yurovsky, D. (2017). A Collaborative Approach to Infant Research: Promoting
522 Reproducibility, Best Practices, and Theory-Building. *Infancy*, 22(4), 421–435.
523 <https://doi.org/10.1111/infa.12182>

524 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank:
525 An open repository for developmental vocabulary data. *Journal of Child
526 Language*, 44(3), 677–694.

527 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability
528 and Consistency in Early Language Learning: The Wordbank Project*. Cambridge,
529 MA: MIT Press.

530 Golinkoff, R. M., Ma, W., Song, L., & Hirsh-Pasek, K. (2013). Twenty-five years
531 using the intermodal preferential looking paradigm to study language acquisition:
532 What have we learned? *Perspectives on Psychol. Science*, 8(3), 316–339.

533 Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P.,
534 ... Poldrack, R. A. (2016). The brain imaging data structure, a format for
535 organizing and describing outputs of neuroimaging experiments. *Scientific Data*,
536 3(1), 160044. <https://doi.org/10.1038/sdata.2016.44>

- 537 Hardwicke, T. E., Mathur, M. B., MacDonald, K., Nilsonne, G., Banks, G. C.,
538 Kidwell, M. C., ... Frank, M. C. (2018). Data availability, reusability, and
539 analytic reproducibility: Evaluating the impact of a mandatory open data policy
540 at the journal *Cognition*. *Royal Society Open Science*, 5(8).
541 <https://doi.org/10.1098/rsos.180448>
- 542 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world?
543 *The Behavioral and Brain Sciences*, 33(2-3), 61–83.
544 <https://doi.org/10.1017/S0140525X0999152X>
- 545 Hirsh-Pasek, K., Cauley, K. M., Golinkoff, R. M., & Gordon, L. (1987). The eyes
546 have it: Lexical and syntactic comprehension in a new paradigm. *Journal of Child
547 Language*, 14(1), 23–45.
- 548 Hurtado, N., Marchman, V. A., & Fernald, A. (2007). Spoken word recognition by
549 Latino children learning Spanish as their first language. *Journal of Child
550 Language*, 34(2), 227–249. <https://doi.org/10.1017/S0305000906007896>
- 551 Hurtado, N., Marchman, V. A., & Fernald, A. (2008). Does input influence uptake?
552 Links between maternal talk, processing speed and vocabulary size in
553 Spanish-learning children. *Developmental Science*, 11(6), 31–39.
554 <https://doi.org/10.1111/j.1467-7687.2008.00768.x>
- 555 Lewis, M., Braginsky, M., Tsuji, S., Bergmann, C., Piccinini, P. E., Cristia, A., &
556 Frank, M. C. (2016). *A Quantitative Synthesis of Early Language Acquisition
557 Using Meta-Analysis* (pp. 1–24). <https://doi.org/10.31234/osf.io/htsjm>
- 558 Lew-Williams, C., & Fernald, A. (2007). Young children learning Spanish make rapid
559 use of grammatical gender in spoken word recognition. *Psychological Science*,
560 18(3), 193–198.

- 561 Mahr, T., McMillan, B. T. M., Saffran, J. R., Ellis Weismer, S., & Edwards, J. (2015).
562 Anticipatory coarticulation facilitates word recognition in toddlers. *Cognition*,
563 142, 345–350. <https://doi.org/10.1016/j.cognition.2015.05.009>
- 564 Marchman, V. A., Loi, E. C., Adams, K. A., Ashland, M., Fernald, A., & Feldman, H.
565 M. (2018). Speed of language comprehension at 18 months old predicts
566 school-relevant outcomes at 54 months old in children born preterm. *Journal of*
567 *Dev. & Behav. Pediatrics*, 39(3), 246–253.
- 568 Muthukrishna, M., Bell, A. V., Henrich, J., Curtin, C. M., Gedranovich, A.,
569 McInerney, J., & Thue, B. (2020). Beyond Western, Educated, Industrial, Rich,
570 and Democratic (WEIRD) Psychology: Measuring and Mapping Scales of
571 Cultural and Psychological Distance. *Psychological Science*, 31(6), 678–701.
- 572 Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A.,
573 ... Vazire, S. (2021). Replicability, Robustness, and Reproducibility in
574 Psychological Science. *PsyArXiv*.
575 <https://doi.org/https://doi.org/10.31234/osf.io/ksfvq>
- 576 Peter, M. S., Durrant, S., Jessop, A., Bidgood, A., Pine, J. M., & Rowland, C. F.
577 (2019). Does speed of processing or vocabulary size predict later language growth
578 in toddlers? *Cognitive Psychology*, 115, 101238.
- 579 R Core Team. (2020). *R: A language and environment for statistical computing*.
580 Vienna, Austria: R Foundation for Statistical Computing. Retrieved from
581 <https://www.R-project.org/>
- 582 Ronfard, S., Wei, R., & Rowe, M. L. (2021). Exploring the linguistic, cognitive, and
583 social skills underlying lexical processing efficiency as measured by the
584 looking-while-listening paradigm. *Journal of Child Language*, 1–24.

- 585 https://doi.org/10.1017/S0305000921000106
- 586 Sanchez, A., Meylan, S. C., Braginsky, M., MacDonald, K. E., Yurovsky, D., & Frank,
587 M. C. (2019). childe-db: A flexible and reproducible interface to the child
588 language data exchange system. *Behavior Research Methods*, 51(4), 1928–1941.
589 https://doi.org/10.3758/s13428-018-1176-7
- 590 Swingley, D., & Aslin, R. N. (2002). Lexical neighborhoods and the word-form
591 representations of 14-month-olds. *Psychological Science*, 13(5), 480–484.
592 https://doi.org/10.1111/1467-9280.00485
- 593 The ManyBabies Consortium. (2020). Quantifying sources of variability in infancy
594 research using the infant-directed speech preference. *Advances in Methods and*
595 *Practices in Psychological Science*, 3(1), 24–52.
- 596 Weisleder, A., & Fernald, A. (2013). Talking to Children Matters: Early Language
597 Experience Strengthens Processing and Builds Vocabulary. *Psychological Science*,
598 24(11), 2143–2152. https://doi.org/10.1177/0956797613488145
- 599 Zettersten, M., Bergey, C., Bhatt, N., Boyce, V., Braginsky, M., Carstensen, A., ...
600 others. (2021). *Peekbank: Exploring children's word recognition through an open,*
601 *large-scale repository for developmental eye-tracking data*.