

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

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

Abstract

20 The ability to rapidly recognize words and link them to referents in context is central to
21 children's early language development. This ability, often called word recognition in the
22 developmental literature, is typically studied in the looking-while-listening paradigm, which
23 measures infants' fixation on a target object (vs. a distractor) after hearing a target label.
24 We present a large-scale, open database of infant and toddler eye-tracking data from
25 looking-while-listening tasks. The goal of this effort is to address theoretical and
26 methodological challenges in measuring vocabulary development.

27 *Keywords:* tools; processing; analysis / usage examples

28 Word count: X

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

31 Across their first years of life, children learn words at an accelerating pace (Frank,
32 Braginsky, Yurovsky, & Marchman, 2021). Although many children will only produce their
33 first word at around one year of age, they show signs of understanding many common
34 nouns (e.g., "mommy") and phrases (e.g., "Let's go bye-bye!") much earlier in development
35 (Bergelson & Swingley, 2012). However, the processes involved in early word understanding
36 are less directly apparent in children's behaviors and are less accessible to observation than
37 developments in speech production (Fernald, Zangl, Portillo, & Marchman, 2008). To
38 understand speech, children must process the incoming auditory signal and link that signal
39 to relevant meanings – a process often referred to as word recognition. Measuring early
40 word recognition offers insight into children's early word representations and as well as the
41 speed and efficiency with which children comprehend language in real time, as the speech
42 signal unfolds (Bergelson, 2020; Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998).
43 Word recognition skills are also thought to build a foundation for children's subsequent
44 language development. Past research has found that early word recognition efficiency is
45 predictive of later linguistic and general cognitive outcomes (Bleses, Makransky, Dale,
46 Højen, & Ari, 2016; Marchman et al., 2018). One explanation for this relationship is that
47 efficiency of word recognition facilitates subsequent word learning: the faster children are
48 at processing speech, the more efficiently they can learn from the input in their
49 environment (Fernald & Marchman, 2012).

50 While word recognition is a central part of children's language development, mapping
51 the trajectory of word recognition skills has remained elusive. Studies investigating
52 children's word recognition are typically limited in scope to experiments in individual labs
53 involving small samples tested on a small set of items. This limitation makes it difficult to
54 understand developmental changes in children's word knowledge at a broad scale.

55 Peekbank provides an openly accessible database of eye-tracking data of children's word
56 recognition, with the primary goal of facilitating the study of developmental changes in
57 children's word knowledge and recognition speed.

58 **The “Looking-While-Listening” Paradigm**

59 Word recognition is traditionally studied in the “looking-while-listening” paradigm
60 [alternatively referred to as the intermodal preferential looking procedure; Fernald, Zangl,
61 Portillo, and Marchman (2008); Hirsh-Pasek, Cauley, Golinkoff, and Gordon (1987)]. In
62 such studies, infants listen to a sentence prompting a specific referent (e.g., “Look at the
63 dog!”) while viewing two images on the screen (e.g., an image of a dog – the target image –
64 and an image of a bird – the distractor image). Infants’ word recognition is measured in
65 terms of how quickly and accurately they fixate on the correct target image after hearing
66 its label. Past research has used this same basic method to study a wide range of questions
67 in language development. For example, the looking-while-listening paradigm has been used
68 to investigate early noun knowledge, phonological representations of words, prediction
69 during language processing, and individual differences in language development (Bergelson
70 & Swingley, 2012; Golinkoff, Ma, Song, & Hirsh-Pasek, 2013; Lew-Williams & Fernald,
71 2007; Marchman et al., 2018; Swingley & Aslin, 2000a).

72 **TO DO: ALIGN CHALLENGES WITH tidybits/ use cases - computational
73 reproducibility and teaching - item-level analyses**

74 **Measuring developmental change in word recognition**

75 While the looking-while-listening paradigm has been highly fruitful in advancing
76 understanding of early word knowledge, fundamental questions remain. One central
77 question is how to accurately capture developmental change in the speed and accuracy of
78 word recognition. There is ample evidence demonstrating that infants get faster and more

79 accurate in word recognition over the first few years of life (e.g., Fernald, Pinto, Swingley,
80 Weinberg, & McRoberts, 1998). However, precisely measuring developmental increases in
81 the speed and accuracy of word recognition remains challenging due to the difficulty of
82 distinguishing developmental changes in word recognition skill from changes in knowledge
83 of specific words. This problem is particularly thorny in studies with young children, since
84 the number of items that can be tested within a single session is limited and items must be
85 selected in an age-appropriate manner (Peter et al., 2019). One way to overcome this
86 challenge is to measure word recognition across development in a large-scale dataset with a
87 wide range of items. A sufficiently large dataset would allow researchers to estimate
88 developmental change in word recognition speed and accuracy while generalizing across
89 changes related to specific words.

90 **Developing methodological best-practices**

91 A second question relates to evaluating methodological best practices. In particular,
92 many fundamental analytic decisions vary substantially across studies, and different
93 decisions may lead to researchers drawing different inferences about children's word
94 recognition. For example, researchers vary in how they select time windows for analysis,
95 transform the dependent measure of target fixations, and model the time course of word
96 recognition (Csibra, Hernik, Mascaro, Tatone, & Lengyel, 2016; Fernald, Zangl, Portillo, &
97 Marchman, 2008; Huang & Snedeker, 2020). This problem is made more complex by the
98 fact that many of these decisions depend on a variety of design-related and
99 participant-related factors (e.g., infant age). Establishing best practices therefore requires a
100 large database of infant word recognition studies varying across such factors, in order to
101 test the potential consequences of methodological decisions on study results.

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

103 What these two questions share is that they are difficult to answer at the scale of a
104 single study. To address this challenge, we introduce Peekbank, a flexible and reproducible
105 interface to an open database of developmental eye-tracking studies. The Peekbank project
106 (a) collects a large set of eye-tracking datasets on children’s word recognition, (b)
107 introduces a data format and processing tools for standardizing eye-tracking data across
108 data sources, and (c) provides an interface for accessing and analyzing the database. In the
109 current paper, we give an overview of the key components of the project and some initial
110 demonstrations of its utility in advancing theoretical and methodological insights. We
111 report two analyses using the database and associated tools ($N=1,233$): (1) a growth curve
112 analysis modeling age-related changes in infants’ word recognition while generalizing across
113 item-level variability; and (2) a multiverse-style analysis of how a central methodological
114 decision – selecting the time window of analysis – impacts inter-item reliability.

115 Design and Technical Approach**116 Database Framework**

117 One of the main challenges in compiling a large-scale eye-tracking dataset is the lack
118 of a shared data format across individual experiments. Researcher conventions for
119 structuring data vary, as do the technical specifications of different devices (e.g., computer
120 displays and eyetracking cameras), rendering the task of integrating datasets from different
121 labs and data sources difficult. Therefore, our first effort was to develop a common tabular
122 format to support analyses of all studies simultaneously.

123 As illustrated in Figure 1, the Peekbank framework consists of four main components:
124 (1) a set of tools to convert eye-tracking datasets into a unified format, (2) a relational
125 database populated with data in this unified format, (3) a set of tools to retrieve data from
126 this database, and (4) a web app (using the Shiny framework) for visualizing the data.

127 These components are supported by three packages. The `peekds` package (for the R
128 language; R Core Team (2020)) helps researchers convert existing datasets to use the
129 standardized format of the database. The `peekbank` module (Python) creates a database
130 with the relational schema and populates it with the standardized datasets produced by
131 `peekds`. The database is implemented in MySQL, an industry standard relational
132 database, which may be accessed by a variety of programming languages, and can be
133 hosted on one machine and accessed by many others over the Internet. The `peekbankr`
134 package (R) provides an application programming interface, or API, that offers high-level
135 abstractions for accessing the tabular data stored in Peekbank. Most users will access data
136 through this final package, in which case the details of data formatting and processing are
137 abstracted away from the user.

138 In the following sections, we will begin by providing the details on the database's
139 organization (or *schema*) and the technical implementation on `peekds`. Users who are
140 primarily interested in accessing the database can skip these details and focus on access
141 through the `peekbankr` API and the web apps.

142 Database Schema

143 The `peekbank` database contains two major types of data: (1) metadata regarding
144 the relevant experiment, participant, and trial, and (2) timecourse looking data, detailing
145 where on the screen a child is looking at a given point in time (Fig. 2).

146 Here, we will give an outline of the tables encoding this data. As is common in
147 relational databases, records of similar types (e.g., participants, trials, experiments, coded
148 looks at each timepoint) are grouped into tables, and records of various types are linked
149 through numeric identifiers.

150 **Metadata.** We encode the metadata available for each study and participation in
151 the database. This metadata can be separated into three parts: (1) subject-level

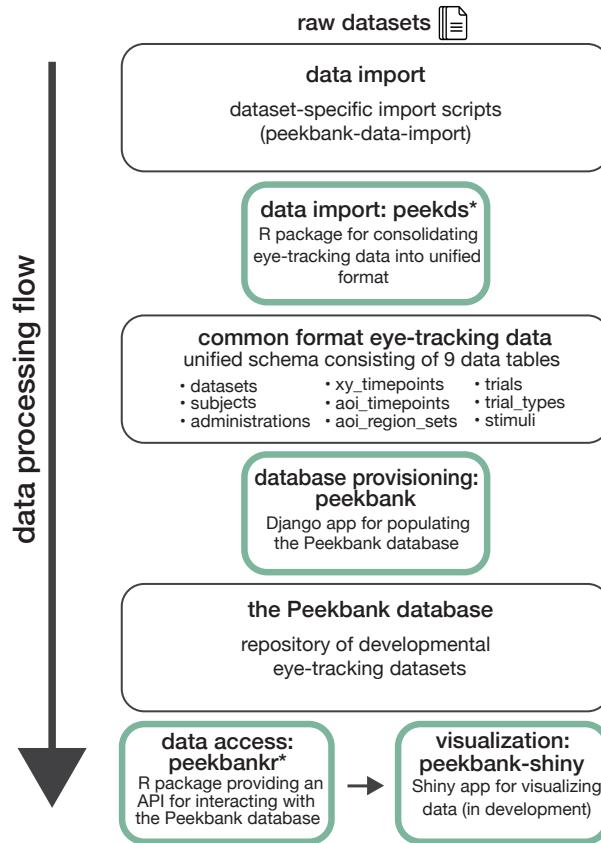


Figure 1. Overview of the Peekbank data ecosystem. Peekbank tools are highlighted in green. * indicates R packages introduced in this work.

152 information (e.g., demographics) (2) experiment-level information (e.g., a subject's age for
 153 a specific experiment, or the particular eyetracker used to collect the data) and (3) trial
 154 information and experimental design (e.g., what images or videos were presented onscreen,
 155 and paired with which audio). Information about individuals who participate in one or
 156 more studies (e.g., a subject's sex and first language), is recorded in the **subjects** table,
 157 while the **administrations** table contains information about a subject's participation in a
 158 single administration of a study (e.g., a subject's age of participation or the eyetracker that
 159 was used). This division allows Peekbank to gracefully handle longitudinal designs: a single
 160 subject can be associated with many administrations.

161 The **stimuli** and **trial_types** tables store information about trials, which in turn
 162 may reflect specifics of the experiment design. Stimuli are (label, image) mappings that are

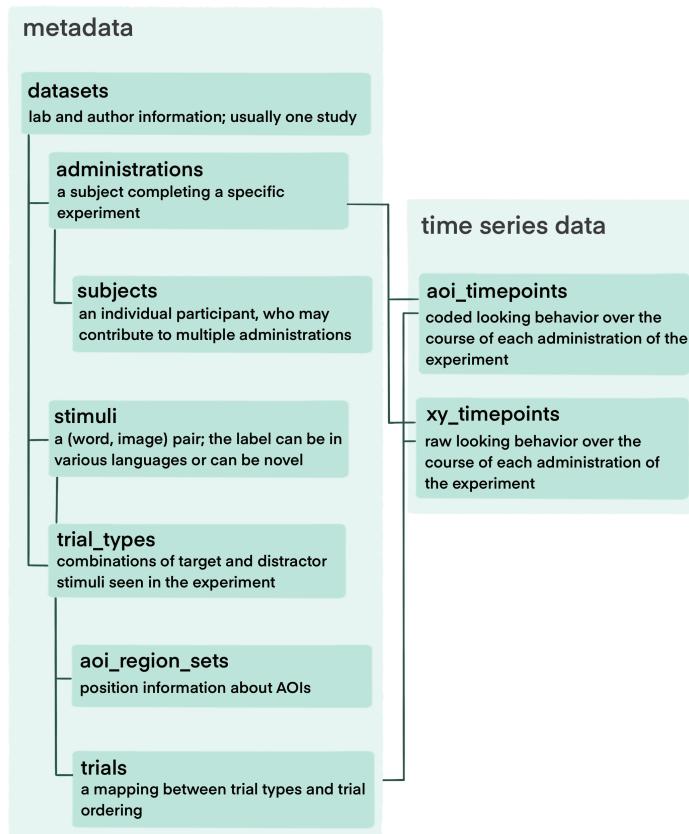


Figure 2. The Peekbank schema. Each square represents a table in the relational database.

seen in the experiment. The `trial_types` table encodes information about each trial of the experiment, including the target stimulus and location, the distractor stimulus and location, and the point of disambiguation for that trial. If this dataset used automatic eyetracking rather than manual coding, each trial type is additionally linked to a set of area of interest (x, y) coordinates, encoded in the `aoi_region_sets` table.

Because individual trial types can be repeated multiple times within an administration, the order of the trials is encoded in the `trials` table. Each unique ordering that occurred in the experiment is encoded in this table. The `trial_id`, which links a trial type to the order it was presented in an administration, is attached to the timecourse looking data.

173 **Timecourse data.** Timecourse looking data is encoded in two tables:

174 `aoi_timepoints` and `xy_timepoints`. The `aoi_timepoints` table encodes where a child is
175 looking at each point in time, by specifying the coded area of interest (AOI): looks to the
176 target, looks to the distractor, looks on the screen but away from target and distractor,
177 and missing looks. All datasets must include this timecourse data, as it represents the
178 main record of children’s looking behavior. For eyetracking experiments that are
179 automatically rather than manually coded, the `xy_timepoints` table additionally encodes
180 the inferred (x, y) coordinates of fixations on the screen over the course of each trial. Both
181 the `aoi_timepoints` and `xy_timepoints` tables are resampled to a consistent sampling
182 rate, as described in the Import section below. To normalize across trials and across
183 experiments, all timecourses are computed so that the time of 0 ms represents the onset of
184 disambiguating material (i.e., the beginning of *dog* in “Can you find the *dog*?”).

185 **Import**

186 During data import, raw eye-tracking datasets are processed to conform to the
187 Peekbank data schema. The following section is a description of the import process for
188 Peekbank. It serves as both a description of our method in importing the datasets already
189 in the database, as well as a high-level overview of the import process for researchers
190 looking to import their data in the future. First, we will describe the import of metadata,
191 and second, we will describe import of the timecourse looking data, including processing
192 functions in `peekds` for normalizing and resampling looking behavior.

193 **Metadata.** Subject-level data is imported for all participants who have experiment
194 data. In general, we import data without particular exclusions, including as many
195 participants as possible in the database. The `subjects` and `administrations` tables
196 separate information at the subject level from information about runs of the experiment,
197 such that longitudinal studies have multiple administrations linked to each subject.

198 The `stimuli` table has a row for each (word, image) pair, and thus is used slightly

199 differently across different experiment designs. In most experiments, there is a one-to-one
200 mapping between images and labels (e.g., each time an image of a dog appears it is referred
201 to as “dog”). For studies in which there are multiple potential labels per image (e.g., “dog”
202 and “chien” are both used to refer to an image of a dog), images can have multiple rows in
203 the `stimuli` table with unique labels as well as a row with no label to be used when the
204 image appears solely as a distractor (and thus its label is ambiguous). This structure is
205 useful for studies on synonymy or using multiple languages. For studies in which the same
206 label refers to multiple images (e.g., the word “dog” refers to an image of a dalmatian and
207 a poodle), the same label can have multiple rows in the `stimuli` table with unique images.
208 The `trial_types` table contains each pair of stimuli, a target and distractor, seen in the
209 experiment. The `trial_types` table links trial types to the `aoi_region_sets` table and
210 the `trials` table.

211 The `trials` table encodes each unique ordering of trial types seen in all runs of an
212 experiment. For example, for experiments with a fixed trial order, the `trials` table will
213 have as many rows as there are stimuli in the experiment; for experiments with a
214 randomized trial order, there will be many rows linking the trial orderings to the trial
215 types. The `trials` table links all experiment design information to the timecourse data.

216 **Timecourse data.** Raw looking data is a series of looks to AOIs or to (x, y)
217 coordinates on the experiment screen, linked to points in time. For data generated by
218 eyetrackers, we typically have (x, y) coordinates at each time point, which will be encoded
219 in the `xy_timepoints` table. These looks will also be recoded into AOIs according to the
220 AOI coordinates in the `aoi_region_sets` table using the `add_aois()` function in `peekds`,
221 which will be encoded in the `aoi_timepoints` table. For hand-coded data, we typically
222 have a series of AOIs; these will be recoded into the categories in the Peekbank schema
223 (target, distractor, other, and missing) and encoded in the `aoi_timepoints` table, and
224 these datasets will not have an `xy_timepoints` table.

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

226 be regularized to center each trial’s time around the point of disambiguation—the time of
227 target word onset in the trial. If time values run throughout the experiment rather than
228 resetting to zero at the beginning of each trial, `rezero_times()` is used to reset the time
229 at each trial. After this, each trial’s times are centered around the point of disambiguation
230 using `normalize_times()`. When these steps are complete, the time course is ready for
231 resampling.

232 To facilitate time course analysis and visualization across datasets, timecourse data
233 must be resampled to a uniform sampling rate (i.e., such that every trial in every dataset
234 has observations at the same time points). To do this, we use the `resample()` function.
235 During the resampling process, we interpolate using constant interpolation, selecting for
236 each interpolated timepoint the looking location for the nearest observed time point in the
237 original data for both `aoi_timepoints` and `xy_timepoints` data. Compared to linear
238 interpolation (see e.g. Wass et al., 2014), constant interpolation has the advantage that it
239 does not introduce new look locations, so it is a more conservative method of resampling.

240 Validation and ingestion into the database

241 After resampling, the final step of dataset import is validation. The `peekds` package
242 offers functions to check the now processed data tables against the database schema to
243 ensure that all tables have the required fields and correct data types for database ingestion.
244 In an effort to double check the data quality and to make sure that no errors are made in
245 the importing script, as part of the import procedure we create a timecourse plot based on
246 our processed tables to replicate the results in the original paper. Once this plot has been
247 created and checked for consistency and all tables pass our validation functions, the
248 processed dataset is ready for ingestion into the database.

249 Currently, the import process is carried out by the Peekbank team using data offered
250 by other research teams. In the future, we hope to allow research teams to carry out their

251 own import processes with checks from the Peekbank team before ingestion. To this end,
 252 import script templates are available for both hand-coded datasets and automatic
 253 eyetracking datasets for research teams to adapt to their data.

254 **CHECK and edit resampling section for ties, interpolating forward/back in**
 255 **time, and for maximum time over which we interpolate**

256 **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
ft_pt	Adams et al., 2018	69	17.1	13–20	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
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
yoursmy	Garrison et al., 2020	35	14.5	12–18	eye-tracking	English

257 The database currently includes 15 looking-while-listening datasets comprising
 258 $N=1320$ total participants (Table 1). Most datasets (12 out of 15 total) consist of data
 259 from monolingual native English speakers. They span a wide age spectrum with
 260 participants ranging from 9 to 70 months of age, and are balanced in terms of gender (46%
 261 female). The datasets vary across a number of design-related dimensions, and include
 262 studies using manually coded video recordings and automated eye-tracking methods (e.g.,
 263 Tobii, EyeLink) to measure gaze behavior. All studies tested familiar items, but the
 264 database also includes 5 datasets that tested novel pseudo-words in addition to familiar
 265 words. All data are openly available on the Open Science Framework
 266 (<https://osf.io/pr6wu/>).

267 How selected? Language coverage? More details about lab and design variation?

268 **Versioning + Expanding the database**

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

276 **Interfacing with peekbank**

277 **Shiny App**

278 One goal of the Peekbank project is to allow a wide range of users to easily explore
279 and learn from the database. We therefore have created an interactive web application –
280 `peekbank-shiny` – that allows users to quickly and easily create informative visualizations
281 of individual datasets and aggregated data. `peekbank-shiny` is built using Shiny, a
282 software package for creating web apps using R. The Shiny app allows users to create
283 commonly used visualizations of looking-while-listening data, based on data from the
284 Peekbank database. Specifically, users can visualize

- 285 1. the time course of looking data in a profile plot depicting infant target looking across
286 trial time
- 287 2. overall accuracy (proportion target looking) within a specified analysis window
- 288 3. reaction times (speed of fixating the target image) in response to a target label
- 289 4. an onset-contingent plot, which shows the time course of participant looking as a
290 function of their look location at the onset of the target label

291 Users are given various customization options for each of these visualizations, e.g.,
292 choosing which datasets to include in the plots, controlling the age range of participants,
293 splitting the visualizations by age bins, and controlling the analysis window for time course
294 analyses. Plots are then updated in real time to reflect users' customization choices, and
295 users are given options to share the visualizations they created. The Shiny app thus allows
296 users to quickly inspect basic properties of Peekbanks datasets and create reproducible
297 visualizations without incurring any of the technical overhead required to access the
298 database through R.

299 **Peekbankr**

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

305 Functions:

- 306 • `connect_to_peekbank()` opens a connection with the Peekbank database to allow
307 tables to be downloaded with the following functions
- 308 • `get_datasets()` gives each dataset name and its citation information
- 309 • `get_subjects()` gives information about persistent subject identifiers (e.g., native
310 languages, sex)
- 311 • `get_administrations()` gives information about specific experimental
312 administrations (e.g., subject age, monitor size, gaze coding method)
- 313 • `get_stimuli()` gives information about word–image pairings that appeared in
314 experiments

- 315 • `get_trial_types()` gives information about pairings of stimuli that appeared in the
- 316 experiment (e.g., point of disambiguation, target and distractor stimuli, condition,
- 317 language)
- 318 • `get_trials()` gives the trial orderings for each administration, linking trial types to
- 319 the trial IDs used in time course data
- 320 • `get_aoi_region_sets()` gives coordinate regions for each area of interest (AOI)
- 321 linked to trial type IDs
- 322 • `get_xy_timepoints()` gives time course data for each subject's looking behavior in
- 323 each trial, as (x, y) coordinates on the experiment monitor
- 324 • `get_aoi_timepoints()` gives time course data for each subject's looking behavior in
- 325 each trial, coded into areas of interest

326 OSF site

327 Stimuli Data in raw format (if some additional datum needed, e.g. pupil size?)

328 **Peekbank: General Descriptives**

329 [Accuracy, Reaction Times, Item variability?]

330 **Overall Word Recognition Accuracy**

Dataset Name	Unique Items	Prop.	Target	95% CI
attword	6	0.63	[0.61, 0.64]	
canine	16	0.64	[0.61, 0.67]	
coartic	10	0.70	[0.67, 0.73]	
cowpig	12	0.60	[0.58, 0.63]	
ft_pt	8	0.64	[0.63, 0.66]	
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	0.62	[0.58, 0.66]	
salientme	16	0.73	[0.71, 0.75]	
switchingCues	40	0.77	[0.75, 0.79]	
tablet	24	0.63	[0.6, 0.67]	
tseatal	30	0.59	[0.54, 0.63]	
yoursmy	87	0.60	[0.56, 0.64]	

Table 2

Average proportion target looking in each dataset.

331 In general, participants demonstrated robust, above-chance word recognition in each
 332 dataset (chance=0.5). Table 2 shows the average proportion of target looking within a
 333 standard critical window of 367-2000ms after the onset of the label for each dataset
 334 (Swingley & Aslin, 2000b). Proportion target looking was generally higher for familiar
 335 words ($M = 0.66$, 95% CI = [0.65, 0.67], $n = 1269$) than for novel words learned during the
 336 experiment ($M = 0.59$, 95% CI = [0.58, 0.61], $n = 822$).

337 **Item-level variability**

338 Figure 3 gives an overview of the variability in accuracy for individual words in each
 339 dataset. The number of unique target labels and their associated accuracy vary widely
 340 across datasets.

341 **Peekbank in Action**

342 We provide two potential use-cases for Peekbank data. In each case, we provide
 343 sample code so as to model how easy it is to do simple analyses using data from the
 344 database. Our first example shows how we can replicate the analysis for a classic study.

345 This type of computational reproducibility can be a very useful exercise for teaching
 346 students about best practices for data analysis (e.g., Hardwicke et al., 2018) and also
 347 provides an easy way to explore looking-while-listening timecourse data in a standardized
 348 format. Our second example shows an in-depth exploration of developmental changes in
 349 the recognition of particular words. Besides its theoretical interest (which we will explore
 350 more fully in subsequent work), this type of analysis could in principle be used for
 351 optimizing the stimuli for new experiments, especially as the Peekbank dataset grows and
 352 gains coverage over a great number of items.

353 Computational reproducibility example: Swingley and Aslin (2000a)

354 Swingley and Aslin (2000a) investigated the specificity of 14-16 month-olds' word
 355 representations using the looking-while-listening paradigm, asking whether recognition
 356 would be slower and less accurate for mispronunciations, e.g. "oppel" (close
 357 mispronunciation) or "opel" (distant mispronunciation) instead of "apple" (correct
 358 condition). In this short vignette, we show how easily the data in Peekbank can be used to
 359 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")
```

360 We begin by retrieving the relevant tables from the database, `aoi_timepoints`,
 361 `administrations`, `trial_types`, and `trials`. As discussed above, each of these can be
 362 downloaded using a simple API call through `peekbankr`, which returns dataframes that
 363 include ID fields. These ID fields allow for easy joining of the data into a single dataframe
 364 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"))

```

As the code above shows, once the data are joined, condition information for each timepoint is present and so we can easily filter out filler trials and set up the conditions for further analysis. For simplicity, here we combine both mispronunciation conditions since this 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()))

```

The final step in our analysis is to create a summary dataframe using `dplyr` commands. We first group the data by timestep, participant, and condition and compute the proportion looking at the correct image. We then summarize again, averaging across participants, computing both means and 95% confidence intervals (via the approximation of 1.96 times the standard error of the mean). The resulting dataframe can be used for visualization of the time-course of looking.

```

ggplot(accuracies, aes(x = t_norm, y = mean_correct, color = condition)) +
  geom_hline(yintercept = 0.5, linetype = "dashed", color = "black") +
  geom_vline(xintercept = 0, linetype = "dotted", color = "black") +
  geom_pointrange(aes(ymin = mean_correct - ci,
                        ymax = mean_correct + ci)) +
  labs(x = "Time from target word onset (msec)",
       y = "Proportion looking at correct image",
       color = "Condition") +
  lims(x = c(-500, 3000))

```

375 Figure 4 shows the average time course of looking for the two conditions, as produced
 376 by the code above. Looks after the correctly pronounced noun appeared both faster
 377 (deviating from chance earlier) and more accurate (showing a higher asymptote). Overall,
 378 this example demonstrates the ability to produce this visualization in just a few lines of
 379 code.

380 **Item analyses**

381 A second use case for Peekbank is to examine item-level variation in word
 382 recognition. Individual datasets rarely have enough statistical power to show reliable
 383 developmental differences within items. To illustrate the power of aggregating data across
 384 multiple datasets, we select the four words with the most data available across studies and
 385 ages (apple, book, dog, and frog) and show average recognition trajectories.

386 Our first step is to collect and join the data from the relevant tables including
 387 timepoint data, trial and stimulus data, and administration data (for participant ages). We
 388 join these into a single dataframe for easy manipulation; this dataframe is a common
 389 starting point for analyses of item-level data.

```
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)

```

390 Next we select a set of four target words (chosen based on having more than XXX
 391 children contributing data for each across several one-year age groups). We create age
 392 groups, aggregate, and compute timepoint-by-timepoint confidence intervals using the z
 393 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())

```

394 Finally, we plot the data as timecourses split by age. Our plotting code is shown
 395 below (with styling commands again removed for clarity). Figure 5 shows the resulting
 396 plot, with time courses for each of three (rather coarse) age bins. Although some baseline
 397 effects are visible across items, we still see clear and consistent increases in looking to the
 398 target, with the increase appearing earlier and in many cases asymptoting at a higher level

399 for older children. On the other hand, this simple averaging approach ignores
 400 study-to-study variation (perhaps responsible for the baseline effects we see in the “apple”
 401 and “frog” items especially). In future work, we hope to introduce model-based analytic
 402 methods that use mixed effects regression to factor out study-level and individual-level
 403 variance in order to recover developmental effects more appropriately (see e.g. Zettersten
 404 et al. (2021) for a prototype of 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)
```

405 Discussion and Conclusion

406 Theoretical progress in understanding child development requires rich datasets, but
 407 collecting child data is expensive, difficult, and time-intensive. Recent years have seen a
 408 growing effort to build open source tools and pool research efforts to meet the challenge of
 409 building a cumulative developmental science (Bergmann et al. (2018); Frank, Braginsky,
 410 Yurovsky, and Marchman (2017); The ManyBabies Consortium (2020)]. The Peekbank
 411 project expands on these efforts by building an infrastructure for aggregating eye-tracking
 412 data across studies, with a specific focus on the looking-while-listening paradigm. This
 413 paper presents an illustration of some of the key theoretical and methodological questions
 414 that can be addressed using Peekbank: generalizing across item-level variability in
 415 children’s word recognition and providing data-driven guidance on methodological choices.

416 There are a number of limitations surrounding the current scope of the database. A
 417 priority in future work will be to expand the size of the database. With 11 datasets

418 currently available in the database, idiosyncrasies of particular designs and condition
419 manipulations still have substantial influence on modeling results. Expanding the set of
420 distinct datasets will allow us to increase the number of observations per item across
421 datasets, leading to more robust generalizations across item-level variability. The current
422 database is also limited by the relatively homogeneous background of its participants, both
423 with respect to language (almost entirely monolingual native English speakers) and
424 cultural background (all but one dataset come from WEIRD populations, potentially
425 limiting generalizability; see Muthukrishna et al. (2020)). Increasing the diversity of
426 participant backgrounds and languages will expand the scope of the generalizations we can
427 form about child word recognition.

428 Finally, while the current database is focused on studies of word recognition, the tools
429 and infrastructure developed in the project can in principle be used to accommodate any
430 eye-tracking paradigm, opening up new avenues for insights into cognitive development.
431 Gaze behavior has been at the core of many of the key advances in our understanding of
432 infant cognition. Aggregating large datasets of infant looking behavior in a single,
433 openly-accessible format promises to bring a fuller picture of infant cognitive development
434 into view.

435 **Acknowledgements**

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

References

- Bergelson, E. (2020). The comprehension boost in early word learning: Older infants are better learners. *Child Development Perspectives*, 14(3), 142–149.
- Bergelson, E., & Swingley, D. (2012). At 6-9 months, human infants know the meanings of many common nouns. *PNAS*, 109(9), 3253–3258.
- Bergmann, C., Tsuji, S., Piccinini, P. E., Lewis, M. L., Braginsky, M., Frank, M. C., & Cristia, A. (2018). Promoting replicability in developmental research through meta-analyses: Insights from language acquisition research. *Child Development*, 89(6), 1996–2009.
- Bleses, D., Makransky, G., Dale, P. S., Højen, A., & Ari, B. A. (2016). Early productive vocabulary predicts academic achievement 10 years later. *Applied Psycholinguistics*, 37(6), 1461–1476.
- Csibra, G., Hernik, M., Mascaro, O., Tatone, D., & Lengyel, M. (2016). Statistical treatment of looking-time data. *Developmental Psychology*, 52(4), 521–536.
- Fernald, A., & Marchman, V. A. (2012). Individual differences in lexical processing at 18 months predict vocabulary growth in typically developing and late-talking toddlers. *Child Development*, 83(1), 203–222.
- Fernald, A., Pinto, J. P., Swingley, D., Weinberg, A., & McRoberts, G. W. (1998). Rapid gains in speed of verbal processing by infants in the 2nd year. *Psychological Science*, 9(3), 228–231.
- Fernald, A., Zangl, R., Portillo, A. L., & Marchman, V. A. (2008). Looking while listening: Using eye movements to monitor spoken language comprehension by infants and young children. In I. A. Sekerina, E. M. Fernandez, & H. Clahsen (Eds.), *Developmental psycholinguistics: On-line methods in children's language processing* (pp. 97–135). Amsterdam: John Benjamins.

- 463 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank:
464 An open repository for developmental vocabulary data. *Journal of Child
465 Language*, 44(3), 677–694.
- 466 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability
467 and Consistency in Early Language Learning: The Wordbank Project*.
468 Cambridge, MA: MIT Press.
- 469 Golinkoff, R. M., Ma, W., Song, L., & Hirsh-Pasek, K. (2013). Twenty-five years
470 using the intermodal preferential looking paradigm to study language acquisition:
471 What have we learned? *Perspectives on Psychol. Science*, 8(3), 316–339.
- 472 Hardwicke, T. E., Mathur, M. B., MacDonald, K., Nilsonne, G., Banks, G. C.,
473 Kidwell, M. C., ... Frank, M. C. (2018). Data availability, reusability, and
474 analytic reproducibility: Evaluating the impact of a mandatory open data policy
475 at the journal Cognition. *Royal Society Open Science*, 5(8).
476 <https://doi.org/10.1098/rsos.180448>
- 477 Hirsh-Pasek, K., Cauley, K. M., Golinkoff, R. M., & Gordon, L. (1987). The eyes
478 have it: Lexical and syntactic comprehension in a new paradigm. *Journal of
479 Child Language*, 14(1), 23–45.
- 480 Huang, Y., & Snedeker, J. (2020). Evidence from the visual world paradigm raises
481 questions about unaccusativity and growth curve analyses. *Cognition*, 200,
482 104251.
- 483 Lew-Williams, C., & Fernald, A. (2007). Young children learning Spanish make
484 rapid use of grammatical gender in spoken word recognition. *Psychological
485 Science*, 18(3), 193–198.
- 486 Marchman, V. A., Loi, E. C., Adams, K. A., Ashland, M., Fernald, A., & Feldman,
487 H. M. (2018). Speed of language comprehension at 18 months old predicts
488 school-relevant outcomes at 54 months old in children born preterm. *Journal of*

- 489 *Dev. & Behav. Pediatrics*, 39(3), 246–253.
- 490 Muthukrishna, M., Bell, A. V., Henrich, J., Curtin, C. M., Gedranovich, A.,
491 McInerney, J., & Thue, B. (2020). Beyond Western, Educated, Industrial, Rich,
492 and Democratic (WEIRD) Psychology: Measuring and Mapping Scales of
493 Cultural and Psychological Distance. *Psychological Science*, 31(6), 678–701.
- 494 Peter, M. S., Durrant, S., Jessop, A., Bidgood, A., Pine, J. M., & Rowland, C. F.
495 (2019). Does speed of processing or vocabulary size predict later language
496 growth in toddlers? *Cognitive Psychology*, 115, 101238.
- 497 R Core Team. (2020). *R: A language and environment for statistical computing*.
498 Vienna, Austria: R Foundation for Statistical Computing. Retrieved from
499 <https://www.R-project.org/>
- 500 Swingley, D., & Aslin, R. N. (2000a). Spoken word recognition and lexical
501 representation in very young children. *Cognition*, 76(2), 147–166.
- 502 Swingley, D., & Aslin, R. N. (2000b). Spoken word recognition and lexical
503 representation in very young children. *Cognition*, 76(2), 147–166.
- 504 The ManyBabies Consortium. (2020). Quantifying sources of variability in infancy
505 research using the infant-directed speech preference. *Advances in Methods and*
506 *Practices in Psychological Science*, 3(1), 24–52.
- 507 Zettersten, M., Bergey, C., Bhatt, N., Boyce, V., Braginsky, M., Carstensen, A., ...
508 others. (2021). Peekbank: Exploring children's word recognition through an
509 open, large-scale repository for developmental eye-tracking data.

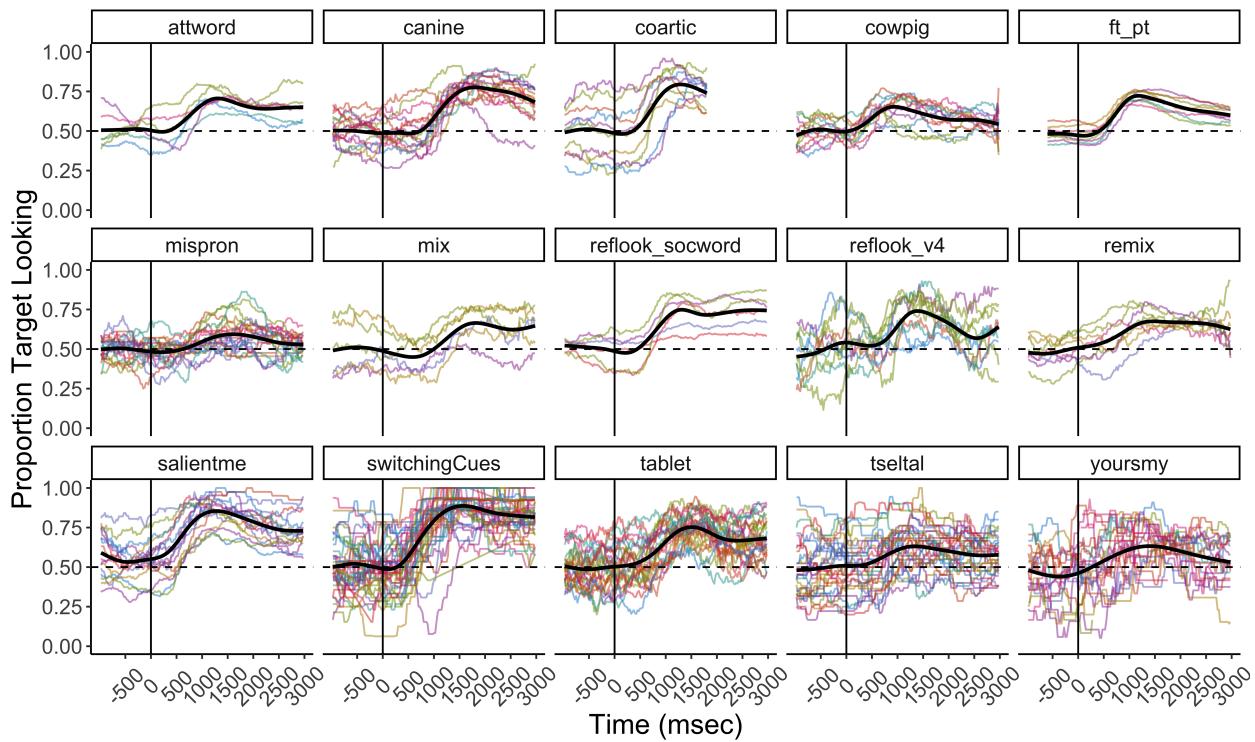


Figure 3. Item-level variability in proportion target looking within each dataset (chance=0.5). Time is centered on the onset of the target label (vertical line). Colored lines represent specific target labels. Black lines represent smoothed average fits based on a general additive model using cubic splines.

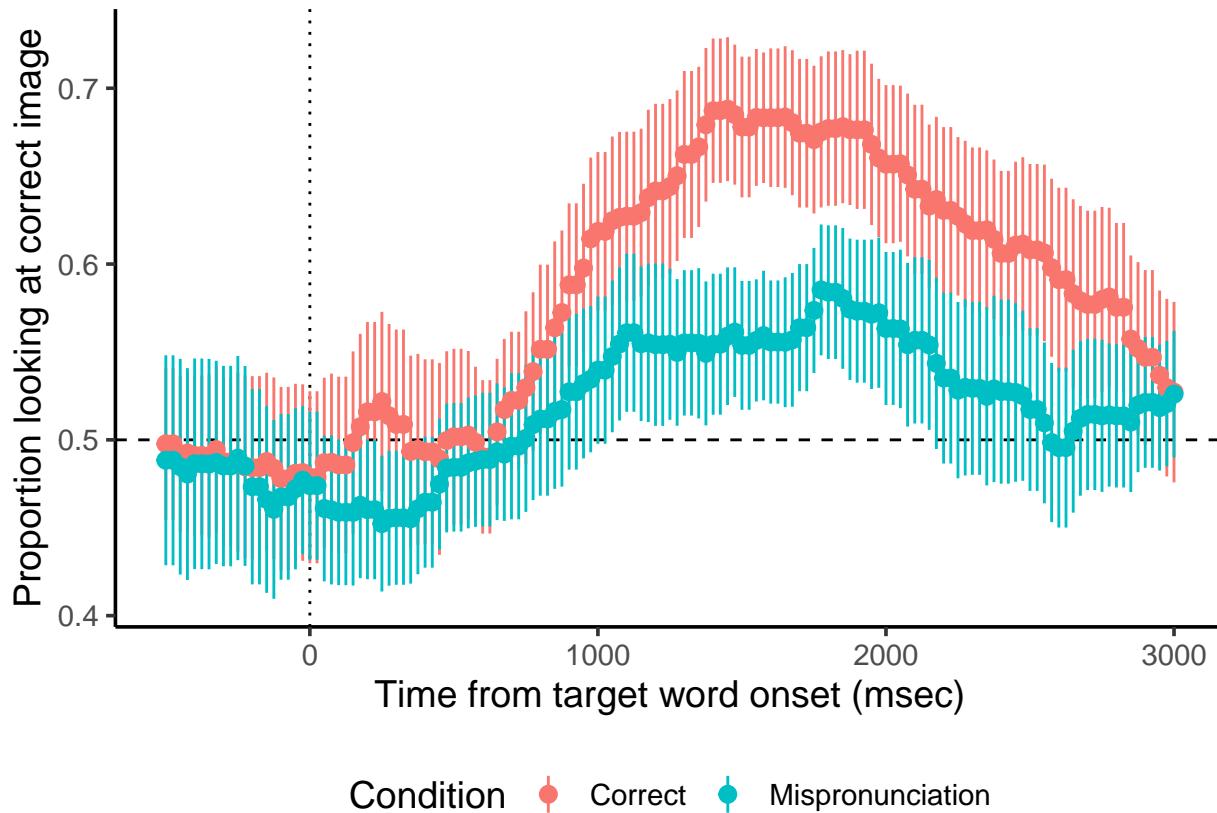


Figure 4. Proportion looking at the correct referent by time from the point of disambiguation (the onset of the target noun). Colors show the two pronunciation conditions; points give means and ranges show 95% confidence intervals. The dotted line shows the point of disambiguation and the dashed line shows chance performance.

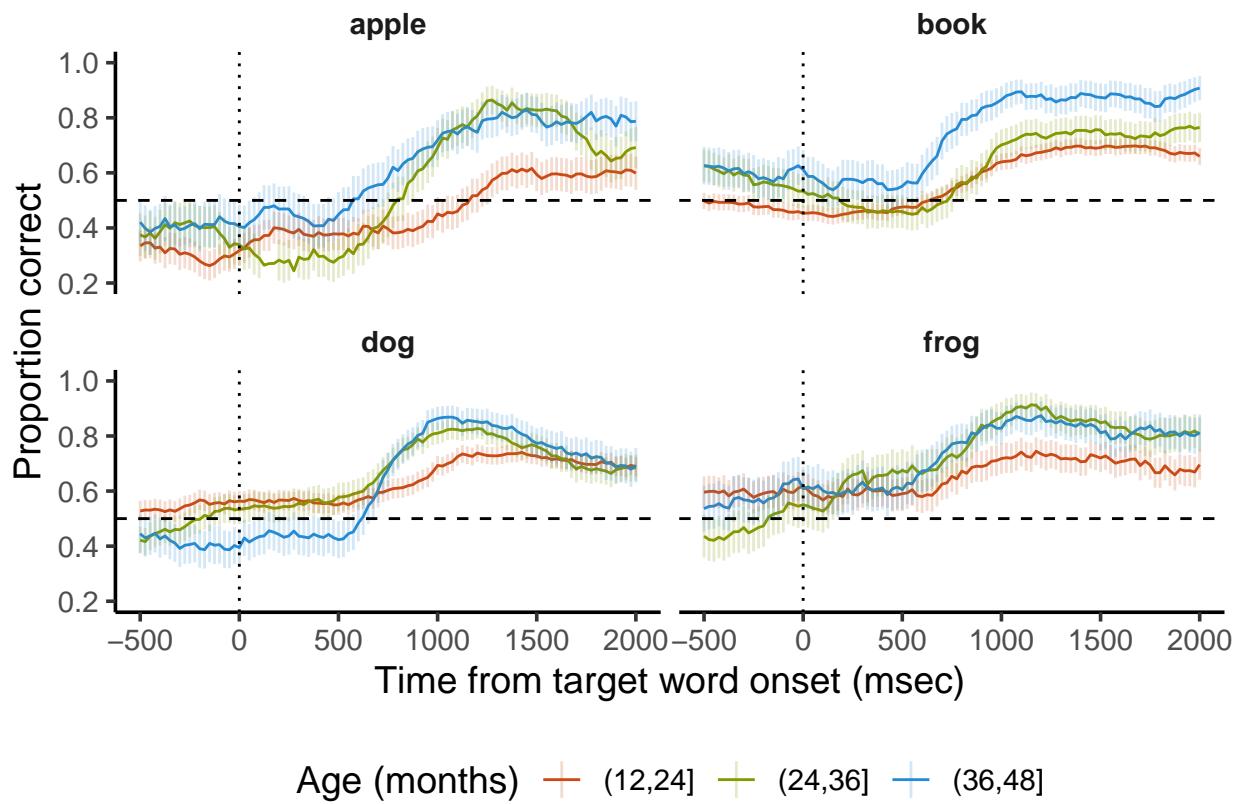


Figure 5. Add caption here.