

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

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 for
30 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 nouns
34 (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 word
40 recognition offers insight into children's early word representations and as well as the speed
41 and efficiency with which children comprehend language in real time, as the speech signal
42 unfolds (Bergelson, 2020; Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998). Word
43 recognition skills are also thought to build a foundation for children's subsequent language
44 development. Past research has found that early word recognition efficiency is predictive of
45 later linguistic and general cognitive outcomes (Bleses, Makransky, Dale, Højen, & Ari, 2016;
46 Marchman et al., 2018).

47 While word recognition is a central part of children's language development, mapping
48 the trajectory of word recognition skills has remained elusive. Studies investigating children's
49 word recognition are typically limited in scope to experiments in individual labs involving
50 small samples tested on a handful of items. The limitations of single datasets makes it
51 difficult to understand developmental changes in children's word knowledge at a broad scale.
52 One way to overcome this challenge is to compile existing datasets into a large-scale
53 database in order to expand the scope of research questions that can be asked about the the
54 development word recognition abilities. This strategy capitalizes on the fact that the

55 looking-while-listening paradigm is widely used, and vast amounts of data have been
56 collected across labs on infants' word recognition over the past 35 years (Golinkoff, Ma, Song,
57 & Hirsh-Pasek, 2013). Such datasets have largely remained isolated from one another, but
58 once combined, they have the potential to offer insights into the lexical development at a
59 broad scale. Similar efforts in language development have born fruit in recent years. For
60 example, WordBank aggregated data from the MacArthur-Bates Communicative
61 Development Inventory, a parent-report measure of child vocabulary, to deliver new insights
62 into cross-linguistic patterns and variability in vocabulary development (Michael C Frank et
63 al., 2017; Frank et al., 2021). In this paper, we introduce *Peekbank*, an open database of
64 infant and toddler eye-tracking data aimed at facilitating the study of developmental
65 changes in children's word knowledge and recognition speed.

66 The “Looking-While-Listening” Paradigm

67 Word recognition is traditionally studied in the “looking-while-listening” paradigm
68 (alternatively referred to as the intermodal preferential looking procedure; Fernald et al.,
69 2008; Hirsh-Pasek, Cauley, Golinkoff, & Gordon, 1987). In such studies, infants listen to a
70 sentence prompting a specific referent (e.g., *Look at the dog!*) while viewing two images on
71 the screen (e.g., an image of a dog – the target image – and an image of a bird – the
72 distractor image). Infants' word recognition is measured in terms of how quickly and
73 accurately they fixate on the correct target image after hearing its label. Past research has
74 used this same basic method to study a wide range of questions in language development.
75 For example, the looking-while-listening paradigm has been used to investigate early noun
76 knowledge, phonological representations of words, prediction during language processing, and
77 individual differences in language development (Bergelson & Swingley, 2012; Golinkoff et al.,
78 2013; Lew-Williams & Fernald, 2007; Marchman et al., 2018; Swingley & Aslin, 2000).

79 Measuring developmental change in word recognition

80 While the looking-while-listening paradigm has been fruitful in advancing
81 understanding of early word knowledge, fundamental questions remain. One central question
82 is how to accurately capture developmental change in the speed and accuracy of word
83 recognition. There is ample evidence demonstrating that infants get faster and more
84 accurate in word recognition over the first few years of life (e.g., Fernald et al., 1998).
85 However, precisely measuring developmental increases in the speed and accuracy of word
86 recognition remains challenging due to the difficulty of distinguishing developmental changes
87 in word recognition skill from changes in knowledge of specific words. This problem is
88 particularly thorny in studies with young children, since the number of items that can be
89 tested within a single session is limited and items must be selected in an age-appropriate
90 manner (Peter et al., 2019). One way to overcome this challenge is to measure word
91 recognition across development in a large-scale dataset with a wide range of items. A
92 sufficiently large dataset would allow researchers to estimate developmental change in word
93 recognition speed and accuracy while generalizing across changes related to specific words.

94 Understanding the development of word recognition on the item level

95 A related open theoretical question is understanding changes in children's word
96 recognition at the level of individual items. Looking-while-listening studies have typically
97 been limited in their ability to assess the development of specific words. One limitation is
98 that studies typically test only a small number of trials for each item, limiting the power the
99 accurately measure the development of word-specific accuracy. A second limitation is that
100 targets are often yoked with a limited set of distractors (often one or two), leaving
101 ambiguous whether accurate looking to a particular target word is largely a function of
102 children's recognition of the target word, their knowledge about the distractor allowing them
103 to reject the distractor as a response candidate, or both. Aggregating across many

104 looking-while-listening studies has the potential to meet these challenges by increasing the
105 number of observations for specific items at different ages and by increasing the variability in
106 the distractor items co-occurring with a specific target.

107 **Replicability and Reproducibility**

108 A core challenge facing psychology in general, and the study of infant development in
109 particular, are threats to the replicability and reproducibility of core empirical results (M. C.
110 Frank et al., 2017; Nosek et al., 2021). In infant research, many studies are not adequately
111 powered to detect the main effects of interest (Bergmann et al., 2018). This is often
112 compounded by low reliability in infant measures, often due to limits on the number of trials
113 that can be collected from an individual infant in an experimental session (Byers-Heinlein,
114 Bergmann, & Savalei, 2021). One hurdle to improving the power in infant research is that it
115 can often be difficult to develop a priori estimates of effect sizes, and how specific design
116 decisions (e.g., the number of test trials) will impact power and reliability. Large-scale
117 databases of infant behavior can aid researchers' in their decision-making by providing rich
118 datasets that can help constrain expectations about possible effect sizes and can be used to
119 make data-driven design decisions. For example, if a researcher is interested in
120 understanding how the number of test trials could impact the power and reliability of their
121 looking-while-listening design, a large-scale database would allow them to simulate possible
122 outcomes across a range of test trials, based on past eye-tracking data with infants.

123 [add paragraph about reproducibility?]

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

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

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

136 **Design and Technical Approach**

137 **Database Framework**

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

144 As illustrated in Figure 1, the Peekbank framework consists of four main components:
145 (1) a set of tools to convert eye-tracking datasets into a unified format, (2) a relational
146 database populated with data in this unified format, (3) a set of tools to retrieve data from
147 this database, and (4) a web app (using the Shiny framework) for visualizing the data. These
148 components are supported by three packages. The `peekds` package (for the R language; R
149 Core Team (2020)) helps researchers convert existing datasets to use the standardized format
150 of the database. The `peekbank` module (Python) creates a database with the relational

151 schema and populates it with the standardized datasets produced by `peekds`. The database
 152 is implemented in MySQL, an industry standard relational database, which may be accessed
 153 by a variety of programming languages, and can be hosted on one machine and accessed by
 154 many others over the Internet. The `peekbankr` package (R) provides an application
 155 programming interface, or API, that offers high-level abstractions for accessing the tabular
 156 data stored in Peekbank. Most users will access data through this final package, in which
 157 case the details of data formatting and processing are abstracted away from the user.

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

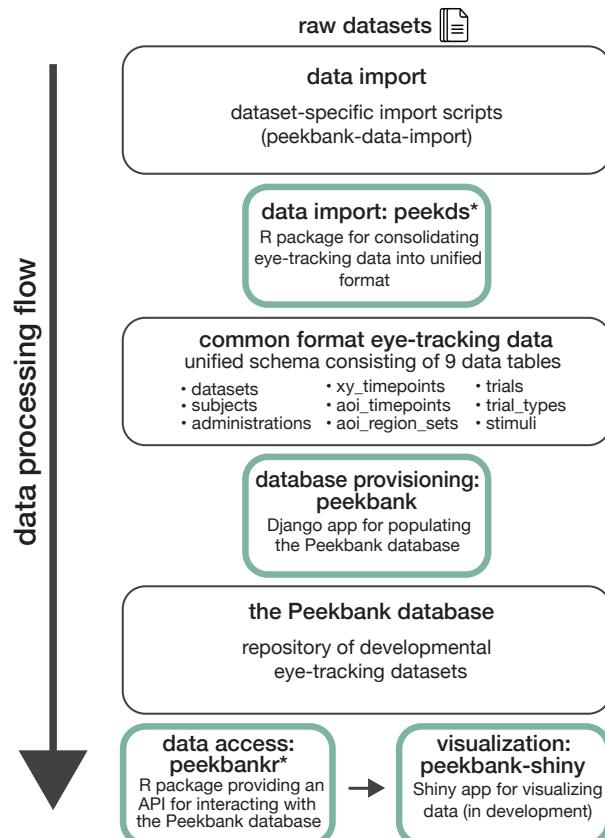


Figure 1. Overview of the Peekbank data ecosystem. Peekbank tools are highlighted in green.
 * indicates R packages introduced in this work.

¹⁶² **Database Schema**

¹⁶³ The peekbank database contains two major types of data: (1) metadata regarding the
¹⁶⁴ relevant experiment, participant, and trial, and (2) time course looking data, detailing where
¹⁶⁵ on the screen a child is looking at a given point in time (Fig. 2).

¹⁶⁶ Here, we will give an outline of the tables encoding this data. As is common in
¹⁶⁷ relational databases, records of similar types (e.g., participants, trials, experiments, coded
¹⁶⁸ looks at each timepoint) are grouped into tables, and records of various types are linked
¹⁶⁹ through numeric identifiers.

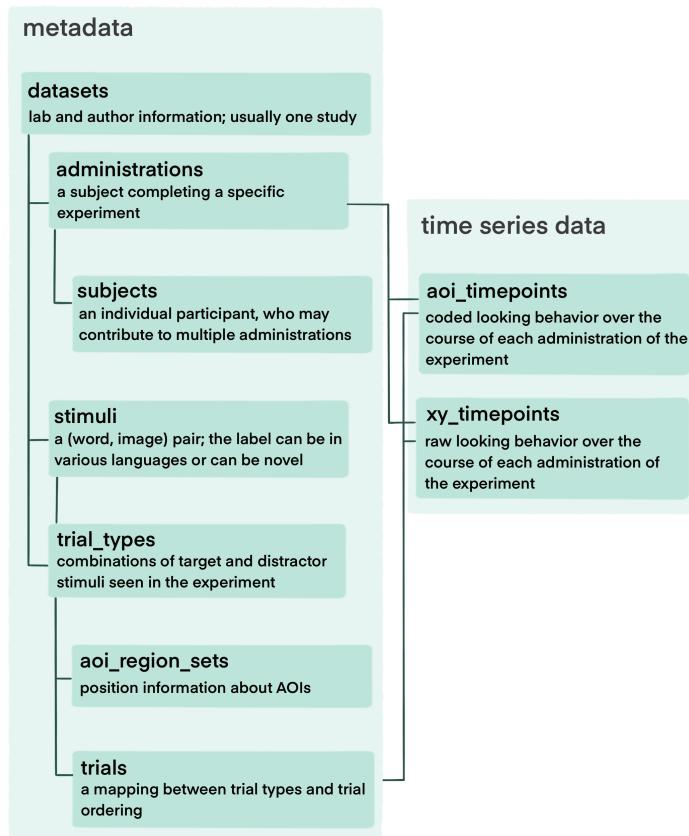


Figure 2. The Peekbank schema. Each square represents a table in the relational database.

¹⁷⁰ **Metadata.** We encode the metadata available for each study and participation in
¹⁷¹ the database. This metadata can be separated into three parts: (1) subject-level information

172 (e.g., demographics) (2) experiment-level information (e.g., a subject's age for a specific
173 experiment, or the particular eye tracker used to collect the data) and (3) trial information
174 and experimental design (e.g., what images or videos were presented onscreen, and paired
175 with which audio). Information about individuals who participate in one or more studies
176 (e.g., a subject's sex and first language), is recorded in the `subjects` table, while the
177 `administrations` table contains information about a subject's participation in a single
178 administration of a study (e.g., a subject's age of participation or the eye tracker that was
179 used). This division allows Peekbank to gracefully handle longitudinal designs: a single
180 subject can be associated with many administrations.

181 The `stimuli` and `trial_types` tables store information about trials, which in turn
182 may reflect specifics of the experiment design. Stimuli are (label, image) mappings that are
183 seen in the experiment. The `trial_types` table encodes information about each trial of the
184 experiment, including the target stimulus and location, the distractor stimulus and location,
185 and the point of disambiguation for that trial. If this dataset used automatic eye-tracking
186 rather than manual coding, each trial type is additionally linked to a set of area of interest
187 (x, y) coordinates, encoded in the `aoi_region_sets` table.

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

192 **Time course data.** Time course looking data is encoded in two tables:
193 `aoi_timepoints` and `xy_timepoints`. The `aoi_timepoints` table encodes where a child is
194 looking at each point in time, by specifying the coded area of interest (AOI): looks to the
195 target, looks to the distractor, looks on the screen but away from target and distractor, and
196 missing looks. All datasets must include this time course data, as it represents the main
197 record of children's looking behavior. For eye-tracking experiments that are automatically

rather than manually coded, the `xy_timepoints` table additionally encodes the inferred (x, y) coordinates of fixations on the screen over the course of each trial. Both the `aoi_timepoints` and `xy_timepoints` tables are resampled to a consistent sampling rate, as described in the Import section below. To normalize across trials and across experiments, all time courses are computed so that the time of 0 ms represents the onset of disambiguating material (i.e., the beginning of *dog* in “Can you find the *dog*?”).

Import

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

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

The `stimuli` table has a row for each (word, image) pair, and thus is used slightly differently across different experiment designs. In most experiments, there is a one-to-one mapping between images and labels (e.g., each time an image of a dog appears it is referred to as “dog”). For studies in which there are multiple potential labels per image (e.g., “dog” and “chien” are both used to refer to an image of a dog), images can have multiple rows in

the `stimuli` table with unique labels as well as a row with no label to be used when the image appears solely as a distractor (and thus its label is ambiguous). This structure is useful for studies on synonymy or using multiple languages. For studies in which the same label refers to multiple images (e.g., the word “dog” refers to an image of a dalmatian and a poodle), the same label can have multiple rows in the `stimuli` table with unique images. The `trial_types` table contains each pair of stimuli, a target and distractor, seen in the experiment. The `trial_types` table links trial types to the `aoi_region_sets` table and the `trials` table.

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

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

Typically, timepoints in the `xy_timepoints` table and `aoi_timepoints` table need to be regularized to center each trial’s time around the point of disambiguation—the time of target word onset in the trial. If time values run throughout the experiment rather than resetting to zero at the beginning of each trial, `rezero_times()` is used to reset the time at

248 each trial. After this, each trial's times are centered around the point of disambiguation
249 using `normalize_times()`. When these steps are complete, the time course is ready for
250 resampling.

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

259 Validation and ingestion into the database

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

268 Currently, the import process is carried out by the Peekbank team using data offered
269 by other research teams. In the future, we hope to allow research teams to carry out their
270 own import processes with checks from the Peekbank team before ingestion. To this end,
271 import script templates are available for both hand-coded datasets and automatic

²⁷² eye-tracking datasets for research teams to adapt to their data.

²⁷³ **CHECK and edit resampling section for ties, interpolating forward/back in
274 time, and for maximum time over which we interpolate**

²⁷⁵ **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

²⁷⁶ The database currently includes 15 looking-while-listening datasets comprising
²⁷⁷ $N=1320$ total participants (Table 1). Most datasets (12 out of 15 total) consist of data from
²⁷⁸ monolingual native English speakers. They span a wide age spectrum with participants
²⁷⁹ ranging from 9 to 70 months of age, and are balanced in terms of gender (46% female). The
²⁸⁰ datasets vary across a number of design-related dimensions, and include studies using
²⁸¹ manually coded video recordings and automated eye-tracking methods (e.g., Tobii, EyeLink)
²⁸² to measure gaze behavior. All studies tested familiar items, but the database also includes 5
²⁸³ datasets that tested novel pseudo-words in addition to familiar words. All data are openly
²⁸⁴ available on the Open Science Framework (<https://osf.io/pr6wu/>).

²⁸⁵ How selected? Language coverage? More details about lab and design variation?

²⁸⁶ **Versioning + Expanding the database**

²⁸⁷ The content of Peekbank will change as we add additional datasets and revise previous
²⁸⁸ ones. To facilitate reproducibility of analyses, we use a versioning system where successive
²⁸⁹ releases are assigned a name reflecting the year and version, e.g., 2021.1. By default, users
²⁹⁰ will interact with the most recent version of the database available, though `peekbankr` API
²⁹¹ allows researchers to run analyses against any previous version of the database. For users
²⁹² with intensive use-cases, each version of the database may be downloaded as a compressed
²⁹³ .sql file and installed on a local MySQL server.

²⁹⁴ **Interfacing with peekbank**

²⁹⁵ **Shiny App**

²⁹⁶ One goal of the Peekbank project is to allow a wide range of users to easily explore and
²⁹⁷ learn from the database. We therefore have created an interactive web application –
²⁹⁸ `peekbank-shiny` – that allows users to quickly and easily create informative visualizations
²⁹⁹ of individual datasets and aggregated data. `peekbank-shiny` is built using Shiny, a software
³⁰⁰ package for creating web apps using R. The Shiny app allows users to create commonly used
³⁰¹ visualizations of looking-while-listening data, based on data from the Peekbank database.
³⁰² Specifically, users can visualize

- ³⁰³ 1. the time course of looking data in a profile plot depicting infant target looking across
³⁰⁴ trial time
- ³⁰⁵ 2. overall accuracy (proportion target looking) within a specified analysis window
- ³⁰⁶ 3. reaction times (speed of fixating the target image) in response to a target label
- ³⁰⁷ 4. an onset-contingent plot, which shows the time course of participant looking as a
³⁰⁸ function of their look location at the onset of the target label

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

317 **Peekbankr**

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

323 Functions:

- 324 • `connect_to_peekbank()` opens a connection with the Peekbank database to allow
325 tables to be downloaded with the following functions
- 326 • `get_datasets()` gives each dataset name and its citation information
- 327 • `get_subjects()` gives information about persistent subject identifiers (e.g., native
328 languages, sex)
- 329 • `get_administrations()` gives information about specific experimental
330 administrations (e.g., subject age, monitor size, gaze coding method)
- 331 • `get_stimuli()` gives information about word–image pairings that appeared in
332 experiments

- 333 • `get_trial_types()` gives information about pairings of stimuli that appeared in the
334 experiment (e.g., point of disambiguation, target and distractor stimuli, condition,
335 language)
- 336 • `get_trials()` gives the trial orderings for each administration, linking trial types to
337 the trial IDs used in time course data
- 338 • `get_aoi_region_sets()` gives coordinate regions for each area of interest (AOI)
339 linked to trial type IDs
- 340 • `get_xy_timepoints()` gives time course data for each subject's looking behavior in
341 each trial, as (x, y) coordinates on the experiment monitor
- 342 • `get_aoi_timepoints()` gives time course data for each subject's looking behavior in
343 each trial, coded into areas of interest

344 OSF site

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

346 Peekbank: General Descriptives

347 [Accuracy, Reaction Times, Item variability?]

³⁴⁸ **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.

³⁴⁹ In general, participants demonstrated robust, above-chance word recognition in each
³⁵⁰ dataset (chance=0.5). Table 2 shows the average proportion of target looking within a
³⁵¹ standard critical window of 367-2000ms after the onset of the label for each dataset (???).
³⁵² Proportion target looking was generally higher for familiar words ($M = 0.66$, 95% CI =
³⁵³ [0.65, 0.67], $n = 1269$) than for novel words learned during the experiment ($M = 0.59$, 95%
³⁵⁴ CI = [0.58, 0.61], $n = 822$).

³⁵⁵ **Item-level variability**

³⁵⁶ Figure 3 gives an overview of the variability in accuracy for individual words in each
³⁵⁷ dataset. The number of unique target labels and their associated accuracy vary widely
³⁵⁸ across datasets.

359

Peekbank in Action

360 We provide two potential use-cases for Peekbank data. In each case, we provide sample
 361 code so as to model how easy it is to do simple analyses using data from the database. Our
 362 first example shows how we can replicate the analysis for a classic study. This type of
 363 computational reproducibility can be a very useful exercise for teaching students about best
 364 practices for data analysis (e.g., Hardwicke et al., 2018) and also provides an easy way to
 365 explore looking-while-listening time course data in a standardized format. Our second
 366 example shows an in-depth exploration of developmental changes in the recognition of
 367 particular words. Besides its theoretical interest (which we will explore more fully in
 368 subsequent work), this type of analysis could in principle be used for optimizing the stimuli
 369 for new experiments, especially as the Peekbank dataset grows and gains coverage over a
 370 greater number of items.

371 Computational reproducibility example: Swingley and Aslin (2000)

372 Swingley and Aslin (2000) investigated the specificity of 14-16 month-olds' word
 373 representations using the looking-while-listening paradigm, asking whether recognition would
 374 be slower and less accurate for mispronunciations, e.g. "oppel" (close mispronunciation) or
 375 "opel" (distant mispronunciation) instead of "apple" (correct pronunciation). In this short
 376 vignette, 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")
```

377 We begin by retrieving the relevant tables from the database, `aoi_timepoints`,
 378 `administrations`, `trial_types`, and `trials`. As discussed above, each of these can be

379 downloaded using a simple API call through `peekbankr`, which returns dataframes that
 380 include ID fields. These ID fields allow for easy joining of the data into a single dataframe
 381 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"))
```

382 As the code above shows, once the data are joined, condition information for each
 383 timepoint is present and so we can easily filter out filler trials and set up the conditions for
 384 further analysis. For simplicity, here we combine both mispronunciation conditions since the
 385 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()))
```

386 The final step in our analysis is to create a summary dataframe using `dplyr`
 387 commands. We first group the data by timestep, participant, and condition and compute the
 388 proportion looking at the correct image. We then summarize again, averaging across
 389 participants, computing both means and 95% confidence intervals (via the approximation of
 390 1.96 times the standard error of the mean). The resulting dataframe can be used for
 391 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))

```

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

Item analyses

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

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

Our first step is to collect and join the data from the relevant tables including timepoint data, trial and stimulus data, and administration data (for participant ages). We join these into a single dataframe for easy manipulation; this dataframe is a common 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)

```

406 Next we select a set of four target words (chosen based on having more than XXX
 407 children contributing data for each across several one-year age groups). We create age
 408 groups, aggregate, and compute timepoint-by-timepoint confidence intervals using the z
 409 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())

```

410 Finally, we plot the data as time courses split by age. Our plotting code is shown

411 below (with styling commands again removed for clarity). Figure 5 shows the resulting plot,
 412 with time courses for each of three (rather coarse) age bins. Although some baseline effects
 413 are visible across items, we still see clear and consistent increases in looking to the target,
 414 with the increase appearing earlier and in many cases asymptoting at a higher level for older
 415 children. On the other hand, this simple averaging approach ignores study-to-study variation
 416 (perhaps responsible for the baseline effects we see in the “apple” and “frog” items
 417 especially). In future work, we hope to introduce model-based analytic methods that use
 418 mixed effects regression to factor out study-level and individual-level variance in order to
 419 recover developmental effects more appropriately (see e.g. Zettersten et al. (2021) for a
 420 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)
```

421

Discussion and Conclusion

422 Theoretical progress in understanding child development requires rich datasets, but
 423 collecting child data is expensive, difficult, and time-intensive. Recent years have seen a
 424 growing effort to build open source tools and pool research efforts to meet the challenge of
 425 building a cumulative developmental science (Bergmann et al. (2018); Frank et al. (2017);
 426 The ManyBabies Consortium (2020)]. The Peekbank project expands on these efforts by
 427 building an infrastructure for aggregating eye-tracking data across studies, with a specific
 428 focus on the looking-while-listening paradigm. This paper presents an illustration of some of
 429 the key theoretical and methodological questions that can be addressed using Peekbank:

430 generalizing across item-level variability in children’s word recognition and providing
431 data-driven guidance on methodological choices.

432 There are a number of limitations surrounding the current scope of the database. A
433 priority in future work will be to expand the size of the database. With 11 datasets currently
434 available in the database, idiosyncrasies of particular designs and condition manipulations
435 still have substantial influence on modeling results. Expanding the set of distinct datasets
436 will allow us to increase the number of observations per item across datasets, leading to more
437 robust generalizations across item-level variability. The current database is also limited by
438 the relatively homogeneous background of its participants, both with respect to language
439 (almost entirely monolingual native English speakers) and cultural background (all but one
440 dataset come from WEIRD populations, potentially limiting generalizability; see
441 Muthukrishna et al. (2020)). Increasing the diversity of participant backgrounds and
442 languages will expand the scope of the generalizations we can form about child word
443 recognition.

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

450 Acknowledgements

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

453

References

- 454 Bergelson, E. (2020). The comprehension boost in early word learning: Older infants are
455 better learners. *Child Development Perspectives*, 14(3), 142–149.
- 456 Bergelson, E., & Swope, D. (2012). At 6-9 months, human infants know the meanings of
457 many common nouns. *PNAS*, 109(9), 3253–3258.
- 458 Bergmann, C., Tsuji, S., Piccinini, P. E., Lewis, M. L., Braginsky, M., Frank, M. C., &
459 Cristia, A. (2018). Promoting replicability in developmental research through
460 meta-analyses: Insights from language acquisition research. *Child Development*,
461 89(6), 1996–2009.
- 462 Bleses, D., Makransky, G., Dale, P. S., Højen, A., & Ari, B. A. (2016). Early productive
463 vocabulary predicts academic achievement 10 years later. *Applied Psycholinguistics*,
464 37(6), 1461–1476.
- 465 Byers-Heinlein, K., Bergmann, C., & Savalei, V. (2021). Six solutions for more reliable infant
466 research. *PsyArXiv*. <https://doi.org/https://doi.org/10.31234/osf.io/ksfvq>
- 467 Fernald, A., Pinto, J. P., Swope, D., Weinberg, A., & McRoberts, G. W. (1998). Rapid
468 gains in speed of verbal processing by infants in the 2nd year. *Psychological Science*,
469 9(3), 228–231.
- 470 Fernald, A., Zangl, R., Portillo, A. L., & Marchman, V. A. (2008). Looking while listening:
471 Using eye movements to monitor spoken language comprehension by infants and
472 young children. In I. A. Sekerina, E. M. Fernandez, & H. Clahsen (Eds.),
473 *Developmental psycholinguistics: On-line methods in children's language processing*
474 (pp. 97–135). Amsterdam: John Benjamins.
- 475 Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., ...

- 476 Yurovsky, D. (2017). A Collaborative Approach to Infant Research: Promoting
477 Reproducibility, Best Practices, and Theory-Building. *Infancy*, 22(4), 421–435.
478 <https://doi.org/10.1111/infa.12182>
- 479 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank: An open
480 repository for developmental vocabulary data. *Journal of Child Language*, 44(3),
481 677–694.
- 482 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability and
483 Consistency in Early Language Learning: The Wordbank Project*. Cambridge, MA:
484 MIT Press.
- 485 Golinkoff, R. M., Ma, W., Song, L., & Hirsh-Pasek, K. (2013). Twenty-five years using the
486 intermodal preferential looking paradigm to study language acquisition: What have
487 we learned? *Perspectives on Psychol. Science*, 8(3), 316–339.
- 488 Hardwicke, T. E., Mathur, M. B., MacDonald, K., Nilsonne, G., Banks, G. C., Kidwell, M.
489 C., . . . Frank, M. C. (2018). Data availability, reusability, and analytic
490 reproducibility: Evaluating the impact of a mandatory open data policy at the journal
491 Cognition. *Royal Society Open Science*, 5(8). <https://doi.org/10.1098/rsos.180448>
- 492 Hirsh-Pasek, K., Cauley, K. M., Golinkoff, R. M., & Gordon, L. (1987). The eyes have it:
493 Lexical and syntactic comprehension in a new paradigm. *Journal of Child Language*,
494 14(1), 23–45.
- 495 Lew-Williams, C., & Fernald, A. (2007). Young children learning Spanish make rapid use of
496 grammatical gender in spoken word recognition. *Psychological Science*, 18(3),
497 193–198.
- 498 Marchman, V. A., Loi, E. C., Adams, K. A., Ashland, M., Fernald, A., & Feldman, H. M.
499 (2018). Speed of language comprehension at 18 months old predicts school-relevant

- 500 outcomes at 54 months old in children born preterm. *Journal of Dev. & Behav.*
501 *Pediatrics*, 39(3), 246–253.
- 502 Muthukrishna, M., Bell, A. V., Henrich, J., Curtin, C. M., Gedranovich, A., McInerney, J.,
503 & Thue, B. (2020). Beyond Western, Educated, Industrial, Rich, and Democratic
504 (WEIRD) Psychology: Measuring and Mapping Scales of Cultural and Psychological
505 Distance. *Psychological Science*, 31(6), 678–701.
- 506 Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A., ...
507 Vazire, S. (2021). Replicability, Robustness, and Reproducibility in Psychological
508 Science. *PsyArXiv*. <https://doi.org/https://doi.org/10.31234/osf.io/ksfvq>
- 509 Peter, M. S., Durrant, S., Jessop, A., Bidgood, A., Pine, J. M., & Rowland, C. F. (2019).
510 Does speed of processing or vocabulary size predict later language growth in toddlers?
511 *Cognitive Psychology*, 115, 101238.
- 512 R Core Team. (2020). *R: A language and environment for statistical computing*. Vienna,
513 Austria: R Foundation for Statistical Computing. Retrieved from
514 <https://www.R-project.org/>
- 515 Swingley, D., & Aslin, R. N. (2000). Spoken word recognition and lexical representation in
516 very young children. *Cognition*, 76(2), 147–166.
- 517 The ManyBabies Consortium. (2020). Quantifying sources of variability in infancy research
518 using the infant-directed speech preference. *Advances in Methods and Practices in*
519 *Psychological Science*, 3(1), 24–52.
- 520 Zettersten, M., Bergey, C., Bhatt, N., Boyce, V., Braginsky, M., Carstensen, A., ... others.
521 (2021). Peekbank: Exploring children's word recognition through an open, large-scale
522 repository for developmental eye-tracking data.

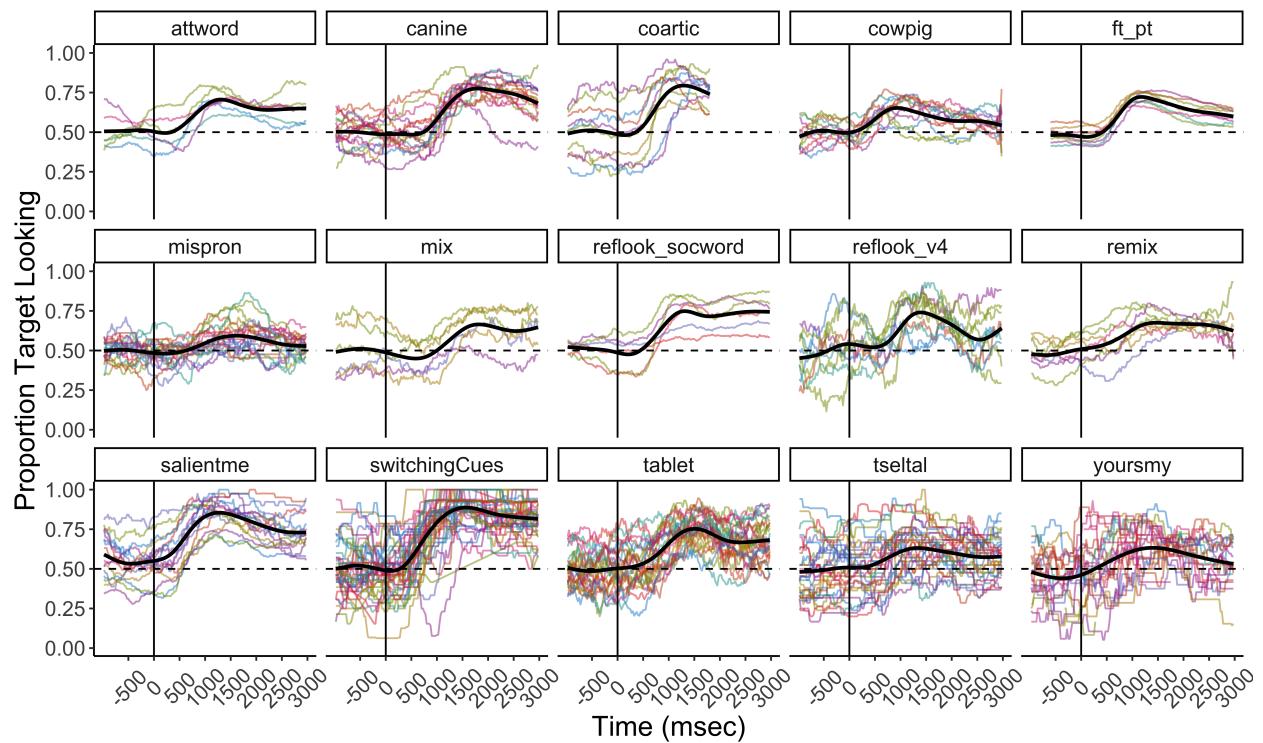


Figure 3. Item-level variability in proportion target looking within each dataset (chance=0.5). Time is centered on the onset of the target label (vertical line). Colored lines represent specific target labels. Black lines represent smoothed average fits based on a general additive model using cubic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

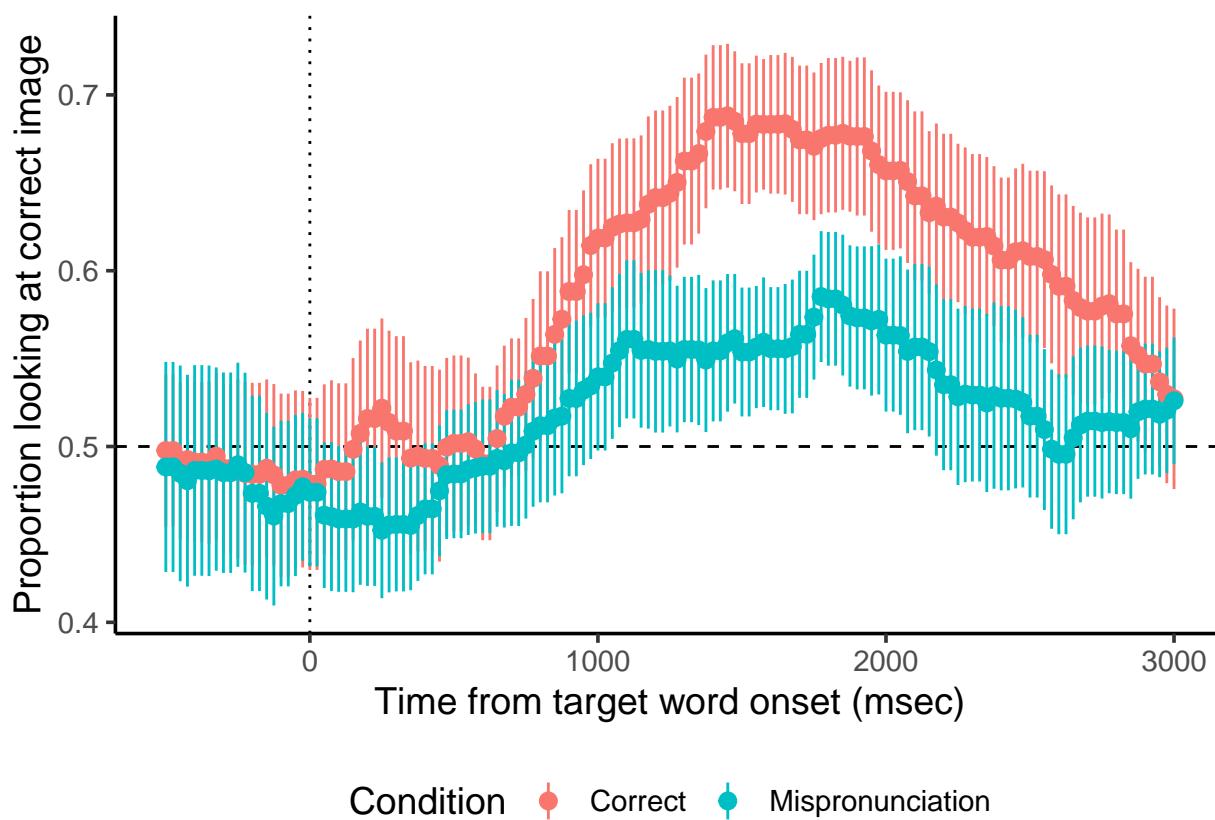


Figure 4. Proportion looking at the correct referent by time from the point of disambiguation (the onset of the target noun). Colors show the two pronunciation conditions; points give means and ranges show 95% confidence intervals. The dotted line shows the point of disambiguation and the dashed line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

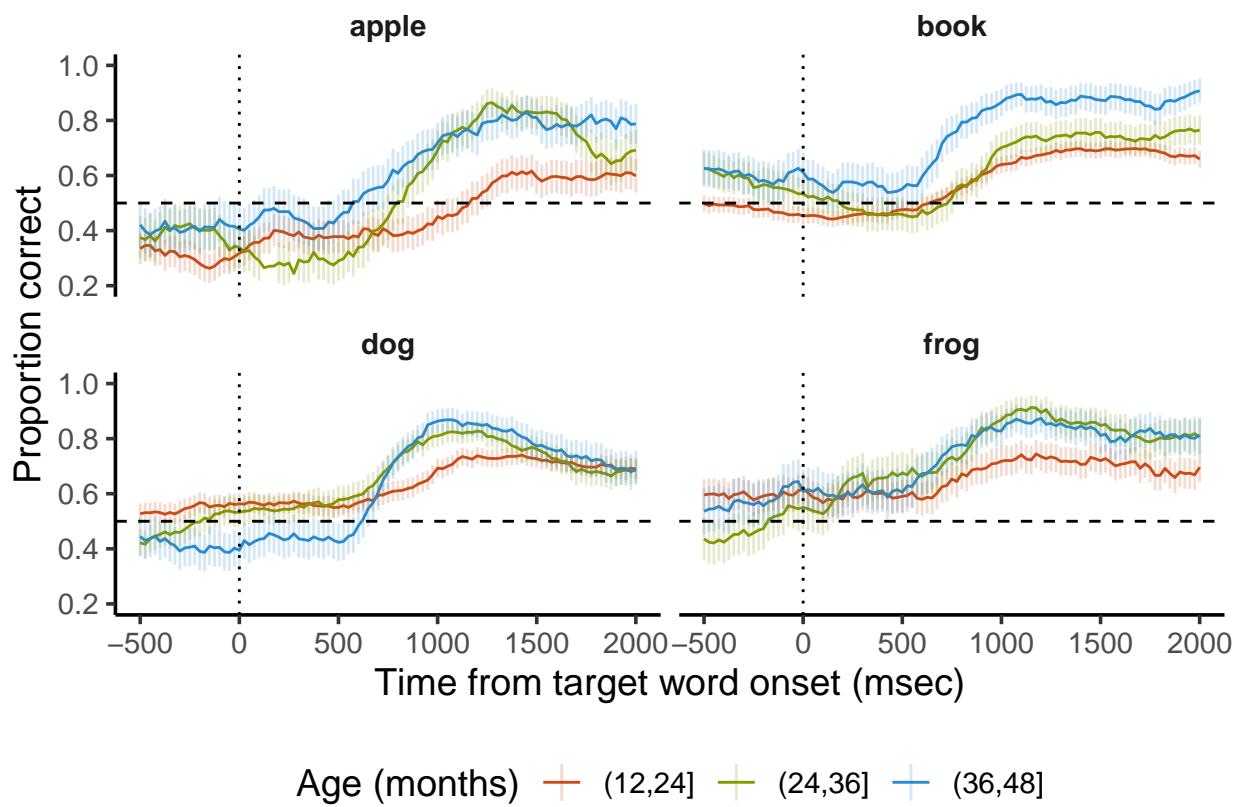


Figure 5. Add caption here.