

<sup>1</sup> Peekbank: An open, large-scale repository for developmental eye-tracking data of children's  
<sup>2</sup> word recognition

<sup>3</sup> Martin Zettersten<sup>1</sup>, Claire Bergey<sup>2</sup>, Naiti S. Bhatt<sup>3</sup>, Veronica Boyce<sup>4</sup>, Mika Braginsky<sup>5</sup>,  
<sup>4</sup> Alexandra Carstensen<sup>4</sup>, Benny deMayo<sup>1</sup>, Kunal Handa<sup>12</sup>, George Kachergis<sup>4</sup>, Molly Lewis<sup>6</sup>,  
<sup>5</sup> Bria Long<sup>4</sup>, Kyle MacDonald<sup>7</sup>, Jessica Mankewitz<sup>4</sup>, Stephan Meylan<sup>5,8</sup>, Annissa N. Saleh<sup>9</sup>,  
<sup>6</sup> Rose M. Schneider<sup>10</sup>, Angeline Sin Mei Tsui<sup>4</sup>, Sarp Uner<sup>8</sup>, Tian Linger Xu<sup>11</sup>, Daniel  
<sup>7</sup> Yurovsky<sup>6</sup>, & Michael C. Frank<sup>4</sup>

<sup>8</sup> <sup>1</sup> Dept. of Psychology, Princeton University

<sup>9</sup> <sup>2</sup> Dept. of Psychology, University of Chicago

<sup>10</sup> <sup>3</sup> Scripps College

<sup>11</sup> <sup>4</sup> Dept. of Psychology, Stanford University

<sup>12</sup> <sup>5</sup> Dept. of Brain and Cognitive Sciences, MIT

<sup>13</sup> <sup>6</sup> Dept. of Psychology, Carnegie Mellon University

<sup>14</sup> <sup>7</sup> Core Technology, McD Tech Labs

<sup>15</sup> <sup>8</sup> Dept. of Psychology and Neuroscience, Duke University

<sup>16</sup> <sup>9</sup> Dept. of Psychology, UT Austin

<sup>17</sup> <sup>10</sup> Dept. of Psychology, UC San Diego

<sup>18</sup> <sup>11</sup> Dept. of Psychological and Brain Sciences, Indiana University

<sup>19</sup> <sup>12</sup> Brown University



20

## Abstract

21 The ability to rapidly recognize words and link them to referents in context is central to  
22 children's early language development. This ability, often called word recognition in the  
23 developmental literature, is typically studied in the looking-while-listening paradigm, which  
24 measures infants' fixation on a target object (vs. a distractor) after hearing a target label.  
25 We present a large-scale, open database of infant and toddler eye-tracking data from  
26 looking-while-listening tasks. The goal of this effort is to address theoretical and  
27 methodological challenges in measuring vocabulary development. We first present how we  
28 created the database, its features and structure, and associated tools for processing and  
29 accessing infant eye-tracking datasets. Using these tools, we then work through two  
30 illustrative examples to show how researchers can use Peekbank to interrogate theoretical  
31 and methodological questions about children's developing word recognition ability.

32       *Keywords:* word recognition; eye-tracking; vocabulary development;  
33 looking-while-listening; visual world paradigm; lexical processing

34 Word count: X

35 Peekbank: An open, large-scale repository for developmental eye-tracking data of children's  
36 word recognition

37 Across their first years of life, children learn words at an accelerating pace (Frank,  
38 Braginsky, Yurovsky, & Marchman, 2021). While many children will only produce their first  
39 word at around one year of age, most children show signs of understanding many common  
40 nouns (e.g., *mommy*) and phrases (e.g., *Let's go bye-bye!*) much earlier in development  
41 (Bergelson & Swingley, 2012, 2013). Although early word understanding is a critical element  
42 of first language learning, the processes involved are less directly apparent in children's  
43 behaviors and are less accessible to observation than developments in speech production  
44 (Fernald, Zangl, Portillo, & Marchman, 2008). To understand a spoken word, children must  
45 process the incoming auditory signal and link that signal to relevant meanings – a process  
46 often referred to as word recognition. One of the primary means of measuring word  
47 recognition in young infants is eye-tracking: gauging where children look in response to  
48 linguistic stimuli to make inferences about children's word processing abilities (Fernald,  
49 Zangl, Portillo, & Marchman, 2008). The logic of this method is that if, upon hearing a  
50 word, a child preferentially looks at a target stimulus rather than a distractor, the child is  
51 able to recognize the word and activate its meaning during real-time language processing.  
52 Measuring early word recognition offers insight into children's early word representations:  
53 children's speed of response (i.e., moving their eyes; turning their heads) to the unfolding  
54 speech signal can reveal children's level of comprehension (Bergelson, 2020; Fernald, Pinto,  
55 Swingley, Weinberg, & McRoberts, 1998). Word recognition skills are also thought to build a  
56 foundation for children's subsequent language development. Past research has found that  
57 early word recognition efficiency is predictive of later linguistic and general cognitive  
58 outcomes (Bleses, Makransky, Dale, Højen, & Ari, 2016; Marchman et al., 2018).

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

61 word recognition are typically limited in scope to experiments in individual labs involving  
62 small samples tested on a handful of items. The limitations of single datasets makes it  
63 difficult to understand developmental changes in children’s word knowledge at a broad scale.

64 One way to overcome this challenge is to compile existing datasets into a large-scale  
65 database in order to expand the scope of research questions that can be asked about the  
66 development of word recognition abilities. This strategy capitalizes on the fact that the  
67 looking-while-listening paradigm is widely used, and vast amounts of data have been  
68 collected across labs on infants’ word recognition over the past 35 years (Golinkoff, Ma, Song,  
69 & Hirsh-Pasek, 2013). Such datasets have largely remained isolated from one another, but  
70 once combined, they have the potential to offer insights into lexical development at a broad  
71 scale. Similar efforts to collect other measures of language development have borne fruit in  
72 recent years. For example, WordBank aggregated data from the MacArthur-Bates  
73 Communicative Development Inventory, a parent-report measure of child vocabulary, to  
74 deliver new insights into cross-linguistic patterns and variability in vocabulary development  
75 (Frank, Braginsky, Yurovsky, & Marchman, 2017, 2021). In this paper, we introduce  
76 *Peekbank*, an open database of infant and toddler eye-tracking data aimed at facilitating the  
77 study of developmental changes in children’s word recognition.

## 78 The “Looking-While-Listening” Paradigm

79 Word recognition is traditionally studied in the “looking-while-listening” paradigm  
80 (Fernald, Zangl, Portillo, & Marchman, 2008; alternatively referred to as the intermodal  
81 preferential looking procedure, Hirsh-Pasek, Cauley, Golinkoff, & Gordon, 1987). In these  
82 studies, infants listen to a sentence prompting a specific referent (e.g., *Look at the dog!*)  
83 while viewing two images on the screen (e.g., an image of a dog – the target image – and an  
84 image of a bird – the distractor image). Infants’ word recognition is evaluated by how  
85 quickly and accurately they fixate on the target image after hearing its label. Past research

86 has used this same basic method to study a wide range of questions in language development.  
87 For example, the looking-while-listening paradigm has been used to investigate early noun  
88 knowledge, phonological representations of words, prediction during language processing, and  
89 individual differences in language development (Bergelson & Swingley, 2012; Golinkoff, Ma,  
90 Song, & Hirsh-Pasek, 2013; Lew-Williams & Fernald, 2007; Marchman et al., 2018; Swingley  
91 & Aslin, 2002).

92 While this research has been fruitful in advancing understanding of early word  
93 knowledge, fundamental questions remain. One central question is how to accurately capture  
94 developmental change in the speed and accuracy of word recognition. There is ample  
95 evidence demonstrating that infants get faster and more accurate in word recognition over  
96 the first few years of life (e.g., Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998).  
97 However, precisely measuring developmental increases in the speed and accuracy of word  
98 recognition remains challenging due to the difficulty of distinguishing developmental changes  
99 in word recognition skill from changes in knowledge of specific words. This problem is  
100 particularly thorny in studies with young children, since the number of items that can be  
101 tested within a single session is limited and items must be selected in an age-appropriate  
102 manner (Peter et al., 2019). More broadly, other key differences in the design choices (e.g.,  
103 how distractor items are selected) and analytic decisions (e.g., how the analysis window is  
104 defined) between studies could obscure developmental change if not appropriately taken into  
105 account.

106 One approach to addressing these challenges is to conduct meta-analyses  
107 aggregating effects across studies while testing for heterogeneity due to researcher choices  
108 (Bergmann et al., 2018; Lewis et al., 2016). However, meta-analyses typically lack the  
109 granularity to estimate participant-level and item-level variation or to model behavior  
110 beyond coarse-grained effect size estimates. An alternative way to approach this challenge is  
111 to aggregate trial-level data from smaller studies measuring word recognition with a wide

range of items and design choices into a large-scale dataset that can be analyzed using a unified modeling approach. A sufficiently large dataset would allow researchers to estimate developmental change in word recognition speed and accuracy while generalizing across changes related to specific words or the design features of particular studies.

A related open theoretical question is understanding changes in children's word recognition at the level of individual items. Looking-while-listening studies have been limited in their ability to assess the development of specific words. One limitation is that studies typically test only a small number of trials for each item, limiting the power to accurately measure the development of word-specific accuracy (DeBolt, Rhemtulla, & Oakes, 2020). A second limitation is that target stimuli are often yoked with a limited set of distractor stimuli (i.e., a child sees a target with only one or two distractor stimuli over the course of an experiment), leaving ambiguous whether accurate looking to a particular target word can be attributed to children's recognition of the target word or their knowledge about the distractor. Aggregating across many looking-while-listening studies has the potential to meet these challenges by increasing the number of observations for specific items at different ages and by increasing the size of the inventory of distractor stimuli that co-occur with each target.

## Replicability and Reproducibility

A core challenge facing psychology in general, and the study of infant development in particular, are threats to the replicability and reproducibility of core empirical results (Frank et al., 2017; Nosek et al., 2021). In infant research, many studies are not adequately powered to detect the main effects of interest (Bergmann et al., 2018). This issue is compounded by low reliability in infant measures, often due to limits on the number of trials that can be collected from an individual infant in an experimental session (Byers-Heinlein, Bergmann, & Savalei, 2021). One hurdle to improving power in infant research is that it can be difficult to

137 develop a priori estimates of effect sizes and how specific design decisions (e.g., the number  
138 of test trials) will impact power and reliability. Large-scale databases of infant behavior can  
139 aid researchers in their decision-making by allowing them to directly test how different  
140 design decisions affect power and reliability. For example, if a researcher is interested in  
141 understanding how the number of test trials could impact the power and reliability of their  
142 looking-while-listening design, a large-scale infant eye-tracking database would allow them to  
143 simulate possible outcomes across a range of test trials, providing the basis for data-driven  
144 design decisions.

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

<sup>164</sup> eye-tracking data.

## <sup>165</sup> **Peekbank: An open database of developmental eye-tracking studies.**

<sup>166</sup> What all of these open challenges share is that they are difficult to address at the scale  
<sup>167</sup> of a single research lab or in a single study. To address this challenge, we developed  
<sup>168</sup> *Peekbank*, a flexible and reproducible interface to an open database of developmental  
<sup>169</sup> eye-tracking studies. The Peekbank project (a) collects a large set of eye-tracking datasets  
<sup>170</sup> on children’s word recognition, (b) introduces a data format and processing tools for  
<sup>171</sup> standardizing eye-tracking data across heterogeneous data sources, and (c) provides an  
<sup>172</sup> interface for accessing and analyzing the database. In the current paper, we introduce the  
<sup>173</sup> key components of the project and give an overview of the existing database. We then  
<sup>174</sup> provide two worked examples of how researchers can use Peekbank. In the first, we examine  
<sup>175</sup> a classic result in the word recognition literature, and in the second we aggregate data across  
<sup>176</sup> studies to investigate developmental trends for the recognition of individual words.

### <sup>177</sup> **Design and Technical Approach**

### <sup>178</sup> **Database Framework**

<sup>179</sup> One of the main challenges in compiling a large-scale eye-tracking database is the lack  
<sup>180</sup> of a shared data format: both labs and individual experiments can record their results in a  
<sup>181</sup> wide range of formats. For example, different experiments encode trial-level and subject-level  
<sup>182</sup> information in many different ways. Therefore, we have developed a common tabular format  
<sup>183</sup> to support analyses of all studies simultaneously.

<sup>184</sup> As illustrated in Figure 1, the Peekbank framework consists of four main components:  
<sup>185</sup> (1) a set of tools to *convert* eye-tracking datasets into a unified format, (2) a relational

186 database populated with data in this unified format, (3) a set of tools to *retrieve* data from  
187 this database, and (4) a web app (using the Shiny framework) for visualizing the data. These  
188 components are supported by three packages. The `peekds` package (for the R language, R  
189 Core Team, 2020) helps researchers convert existing datasets to use the standardized format  
190 of the database. The `peekbank` module (Python) creates a database with the relational  
191 schema and populates it with the standardized datasets produced by `peekds`. The database  
192 is served through MySQL, an industry standard relational database server, which may be  
193 accessed by a variety of programming languages, and can be hosted on one machine and  
194 accessed by many others over the Internet. As is common in relational databases, records of  
195 similar types (e.g., participants, trials, experiments, coded looks at each timepoint) are  
196 grouped into tables, and records of various types are linked through numeric identifiers. The  
197 `peekbankr` package (R) provides an application programming interface, or API, that offers  
198 high-level abstractions for accessing the tabular data stored in Peekbank. Most users will  
199 access data through this final package, in which case the details of data formatting,  
200 processing, and the specifics of connecting to the database are abstracted away from the user.

## 201 Database Schema

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

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

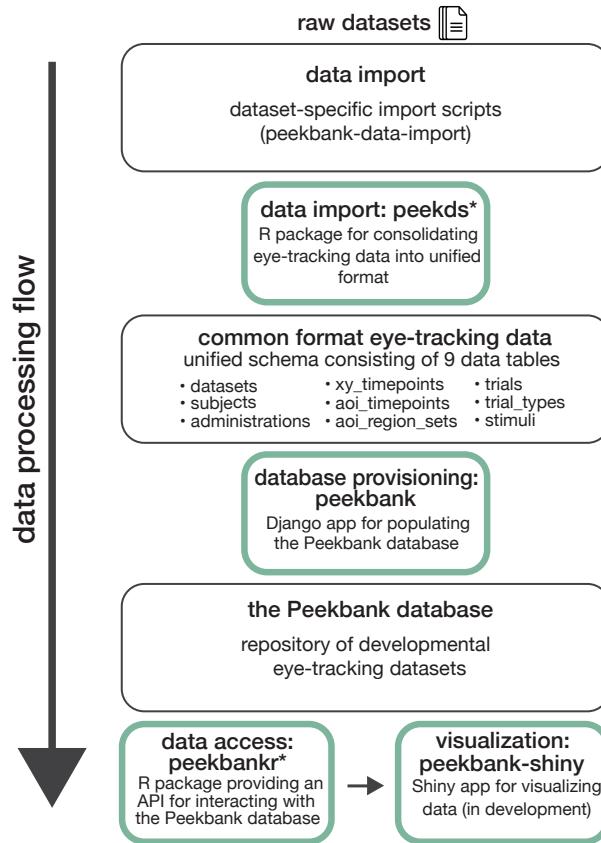


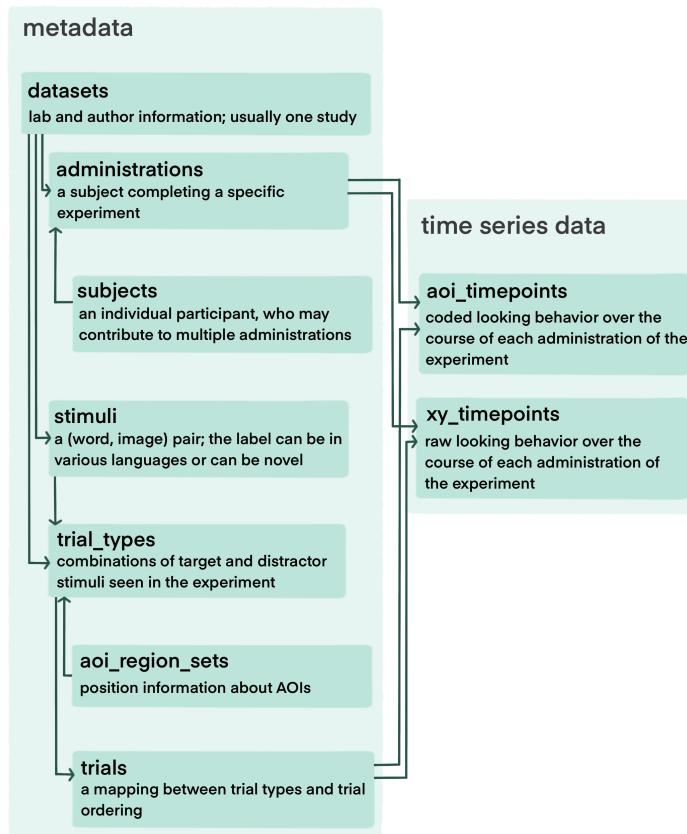
Figure 1. Overview of the Peekbank data ecosystem. Peekbank tools are highlighted in green.  
\* indicates R packages introduced in this work.

### 210 **Participant Information.**

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

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

### 219 **Experiment Information.**



*Figure 2.* The Peekbank schema. Each darker rectangle represents a table in the relational database.

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

223        Information about the experimental design is split across the **trial\_types** and  
 224   **stimuli** tables. The **trial\_types** table encodes information about each trial *in the design*  
 225   *of the experiment*,<sup>\footnote{We note that the term *trial* is ambiguous and could be used to}</sup> refer to both a particular combination of stimuli seen by many participants and a participant  
 226   seeing that particular combination at a particular point in the experiment. We track the  
 227   former in the **trial\_types** table and the latter in the **trials** table.<sup>\footnote{including the target}</sup> including the target  
 228   stimulus and location (left vs. right), the distractor stimulus and location, and the point of  
 229   disambiguation for that trial. If a dataset used automatic eye-tracking rather than manual  
 230

231 coding, each trial type is additionally linked to a set of area of interest (x, y) coordinates,  
232 encoded in the `aoi_region_sets` table. The `trial_types` table links trial types to the  
233 `aoi_region_sets` table and the `trials` table. Each trial\_type record links to two records  
234 in the `stimuli` table, identified by the `distractor_id` and the `target_id` fields.

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

244 ***Session Information.***

245 The `administrations` table includes information about the participant or experiment  
246 that may change between sessions of the same study, even for the same participant. This  
247 includes the age of the participant, the coding method (eye-tracking vs. hand-coding), and  
248 the properties of the monitor that was used.

249 ***Trial Information.***

250 The `trials` table includes information about a specific participant completing a  
251 specific instance of a trial type. This table links each record in the time course looking data  
252 (described below) to the trial type and specifies the order of the trials seen by a specific  
253 participant.

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

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

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

257 coordinates at each time point, which we encode in the `xy_timepoints` table. These looks

258 are also recoded into AOIs according to the AOI coordinates in the `aoi_region_sets` table

259 using the `add_aois()` function in `peekds`, and encoded in the `aoi_timepoints` table. For

260 hand-coded data, we typically have a series of AOIs (i.e., looks to the left vs. right of the

261 screen), but lack information about exact gaze positions on-screen; in these cases the AOIs

262 are recoded into the categories in the Peekbank schema (target, distractor, other, and

263 missing) and encoded in the `aoi_timepoints` table; however, these datasets do not have any

264 corresponding data in the `xy_timepoints` table.

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

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

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

268 *dog?*). We re-centered timing information to the onset of the target label to facilitate

269 comparison of target label processing across all datasets.<sup>1</sup> If time values run throughout the

270 experiment rather than resetting to zero at the beginning of each trial, `rezero_times()` is

271 used to reset the time at each trial. After this, each trial's times are centered around the

272 point of disambiguation using `normalize_times()`. When these steps are complete, the

273 time course is ready for resampling.

274       To facilitate time course analysis and visualization across datasets, time course data

275 must be resampled to a uniform sampling rate (i.e., such that every trial in every dataset has

276 observations at the same time points). All data in the database is resampled to 40 Hz

---

<sup>1</sup> While information preceding the onset of the target label in some datasets such as co-articulation cues (Mahr, McMillan, Saffran, Ellis Weismer, & Edwards, 2015) or adjectives (Fernald, Marchman, & Weisleder, 2013) can in principle disambiguate the target referent, we use a standardized point of disambiguation based on the onset of the label for the target referent. Onset times for other potentially disambiguating information (such as adjectives) can typically be recovered from the raw data provided on OSF.

(observations every 25 ms), which represents a compromise between retaining fine-grained timing information from datasets with dense sampling rates (maximum sampling rate among current datasets: 500 Hz) while minimizing the possibility of introducing artifacts via resampling for datasets with lower sampling rates (minimum sampling rate for current datasets: 30 Hz). Further, 25 ms is a mathematically convenient interval for ensuring consistent resampling; we found that using 33.333 ms (30 Hz) as our interval simply introduced a large number of technical complexities. The resampling operation is accomplished using the `resample_times()` function. During the resampling process, we interpolate using constant interpolation, selecting for each interpolated timepoint the looking location for the earlier-observed time point in the original data for both `aoi_timepoints` and `xy_timepoints` data. Compared to linear interpolation (see e.g., Wass, Smith, & Johnson, 2013), which interpolates between locations of successive timepoints, constant interpolation has the advantage that it is more conservative, in the sense that it does not introduce new (spatial) look locations beyond those measured in the original data. One possible application of our new dataset is investigating the consequences of other interpolation functions for data analysis.

### Processing, Validation, and Ingestion

The `peekds` package offers functions to extract the above data. Once these data have been extracted in a tabular form, the package also offers a function to check whether all tables have the required fields and data types expected by the database. In an effort to double check the data quality and to make sure that no errors are made in the importing script, as part of the import procedure we create a time course plot based on our processed tables to replicate the results in the paper that first presented each dataset. Once this plot has been created and checked for consistency and all tables pass our validation functions, the processed dataset is ready for reprocessing into the database using the `peekbank` library.

<sup>302</sup> This library applies additional data checks, and adds the data to the MySQL database using  
<sup>303</sup> the Django web framework.

<sup>304</sup> Currently, the import process is carried out by the Peekbank team using data offered  
<sup>305</sup> by other research teams. In the future, we hope to allow research teams to carry out their  
<sup>306</sup> own import processes with checks from the Peekbank team before reprocessing. To this end,  
<sup>307</sup> import script templates are available for both hand-coded datasets and automatic  
<sup>308</sup> eye-tracking datasets for research teams to adapt to their data.

## <sup>309</sup> Current Data Sources

Table 1  
*Overview of the datasets in the current database.*

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

<sup>310</sup> The database currently includes 20 looking-while-listening datasets comprising  $N=1594$   
<sup>311</sup> total participants (Table 1). The current data represents a convenience sample of datasets  
<sup>312</sup> that were (a) datasets collected by or available to Peekbank team members, (b) made  
<sup>313</sup> available to Peekbank after informal inquiry or (c) datasets that were openly available. Most  
<sup>314</sup> datasets (14 out of 20 total) consist of data from monolingual native English speakers. They  
<sup>315</sup> span a wide age spectrum with participants ranging from 9 to 70 months of age, and are

316 balanced in terms of gender (47% female). The datasets vary across a number of  
317 design-related dimensions, and include studies using manually coded video recordings and  
318 automated eye-tracking methods (e.g., Tobii, EyeLink) to measure gaze behavior. All studies  
319 tested familiar items, but the database also includes 5 datasets that tested novel  
320 pseudo-words in addition to familiar words. Users interested in a subset of the data (e.g.,  
321 only trials testing familiar words) can filter out unwanted trials using columns available in  
322 the schema (e.g., using the column `stimulus_novelty`).

323 **Versioning and Reproducibility**

324 The content of Peekbank will change as we add additional datasets and revise previous  
325 ones. To facilitate reproducibility of analyses, we use a versioning system by which  
326 successive releases are assigned a name reflecting the year and version, e.g., 2022.1. By  
327 default, users will interact with the most recent version of the database available, though the  
328 `peekbankr` API allows researchers to run analyses against any previous version of the  
329 database. For users with intensive use-cases, each version of the database may be  
330 downloaded as a compressed .sql file and installed on a local MySQL server.

331 Peekbank allows for fully reproducible analyses using our source data, but the goal is  
332 not to reproduce precisely the analyses – or even the datasets – in the publications whose  
333 data we archive. Because of our emphasis on a standardized data importing and formatting  
334 pipeline, there may be minor discrepancies in the timecourse data that we archive compared  
335 with those reported in original publications. Further, we archive all of the data that are  
336 provided to us – including participants that might have been excluded in the original studies,  
337 if these data are available – rather than attempting to reproduce specific exclusion criteria.  
338 We hope that Peekbank can be used as a basis for comparing different exclusion and filtering  
339 criteria – as such, an inclusive policy regarding importing all available data helps us provide  
340 a broad base of data for investigating these decisions.

## 341                   Interfacing with Peekbank

### 342           Peekbankr

343           The `peekbankr` API offers a way for users to access data from the database and  
344           flexibly analyze it in R. The majority of API calls simply allow users to download tables (or  
345           subsets of tables) from the database. In particular, the package offers the following functions:

- 346           • `connect_to_peekbank()` opens a connection with the Peekbank database to allow  
347           tables to be downloaded with the following functions
- 348           • `get_datasets()` gives each dataset name and its citation information
- 349           • `get_subjects()` gives information about persistent subject identifiers (e.g., native  
350           languages, sex)
- 351           • `get_administrations()` gives information about specific experimental  
352           administrations (e.g., subject age, monitor size, gaze coding method)
- 353           • `get_stimuli()` gives information about word–image pairings that appeared in  
354           experiments
- 355           • `get_trial_types()` gives information about pairings of stimuli that appeared in the  
356           experiment (e.g., point of disambiguation, target and distractor stimuli, condition,  
357           language)
- 358           • `get_trials()` gives the trial orderings for each administration, linking trial types to  
359           the trial IDs used in time course data
- 360           • `get_aoi_region_sets()` gives coordinate regions for each area of interest (AOI)  
361           linked to trial type IDs
- 362           • `get_xy_timepoints()` gives time course data for each subject’s looking behavior in  
363           each trial, as (x, y) coordinates on the experiment monitor
- 364           • `get_aoi_timepoints()` gives time course data for each subject’s looking behavior in  
365           each trial, coded into areas of interest

366        Once users have downloaded tables, they can be merged using `join` command via their  
367        linked IDs. A set of standard merges are shown below in the “Peekbank in Action” section;  
368        these allow the common use-case of examining time course data and metadata jointly.

369        Because of the size of the XY and AOI data tables, downloading data across multiple  
370        studies can be time-consuming. Many of the most common analyses of the Peekbank data  
371        require download of the `aoi_timepoints` table, thus we have put substantial work into  
372        optimizing transfer times. In particular, `connect_to_peekbank` offers a data compression  
373        option, and `get_aoi_timepoints` by default downloads time-courses via a compressed  
374        (run-length encoded) representation, which is then uncompressed on the client side. More  
375        information about these options (including how to modify them) can be found in the  
376        package documentation.

## 377        Shiny App

378        One goal of the Peekbank project is to allow a wide range of users to easily explore and  
379        learn from the database. We therefore have created an interactive web application –  
380        `peekbank-shiny` – that allows users to quickly and easily create informative visualizations  
381        of individual datasets and aggregated data. `peekbank-shiny` is built using Shiny, a software  
382        package for creating web apps for data exploration with R, as well as the `peekbankr` package.  
383        The Shiny app allows users to create commonly used visualizations of looking-while-listening  
384        data, based on data from the Peekbank database. Specifically, users can visualize:

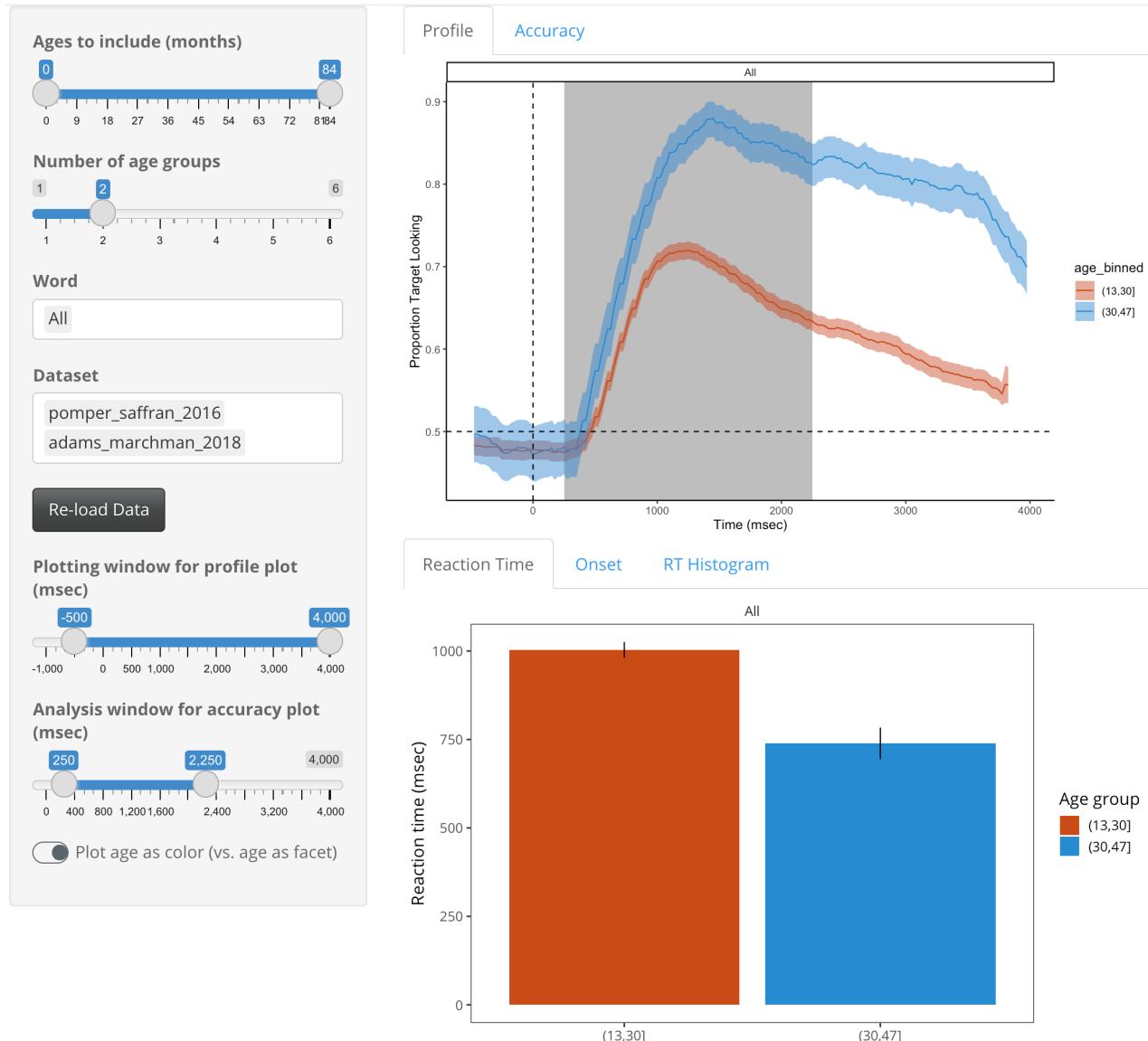
- 385        1. the time course of looking data in a profile plot depicting infant target looking across  
386        trial time
- 387        2. overall accuracy (proportion target looking) within a specified analysis window
- 388        3. reaction times (speed of fixating the target image) in response to a target label
- 389        4. an onset-contingent plot, which shows the time course of participant looking as a

390 function of their look location at the onset of the target label

391 Users are given various customization options for each of these visualizations, e.g.,  
392 choosing which datasets to include in the plots, controlling the age range of participants,  
393 splitting the visualizations by age bins, and controlling the analysis window for time course  
394 analyses. Plots are then updated in real time to reflect users' customization choices, and  
395 users are given options to share the visualizations they created. A screenshot of the app is  
396 shown in Figure 3. The Shiny app thus allows users to quickly inspect basic properties of  
397 Peekbanks datasets and create reproducible visualizations without incurring any of the  
398 technical overhead required to access the database through R.

## 399 OSF site

400 In addition to the Peekbank database proper, all data is openly available on the  
401 Peekbank OSF webpage (<https://osf.io/pr6wu/>). The OSF site also includes the original raw  
402 data (both time series data and metadata, such as trial lists and participant logs) that was  
403 obtained for each study and subsequently processed into the standardized Peekbank format.  
404 Users who are interested in inspecting or reproducing the processing pipeline for a given  
405 dataset can use the respective import script (openly available on GitHub,  
406 <https://github.com/langcog/peekbank-data-import>) to download and process the raw data  
407 from OSF into its final standardized format. Where available, the OSF page also includes  
408 additional information about the stimuli used in each dataset, including in some instances  
409 the original stimulus sets (e.g., image and audio files).



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

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

Table 2  
*Average proportion target looking in each dataset.*

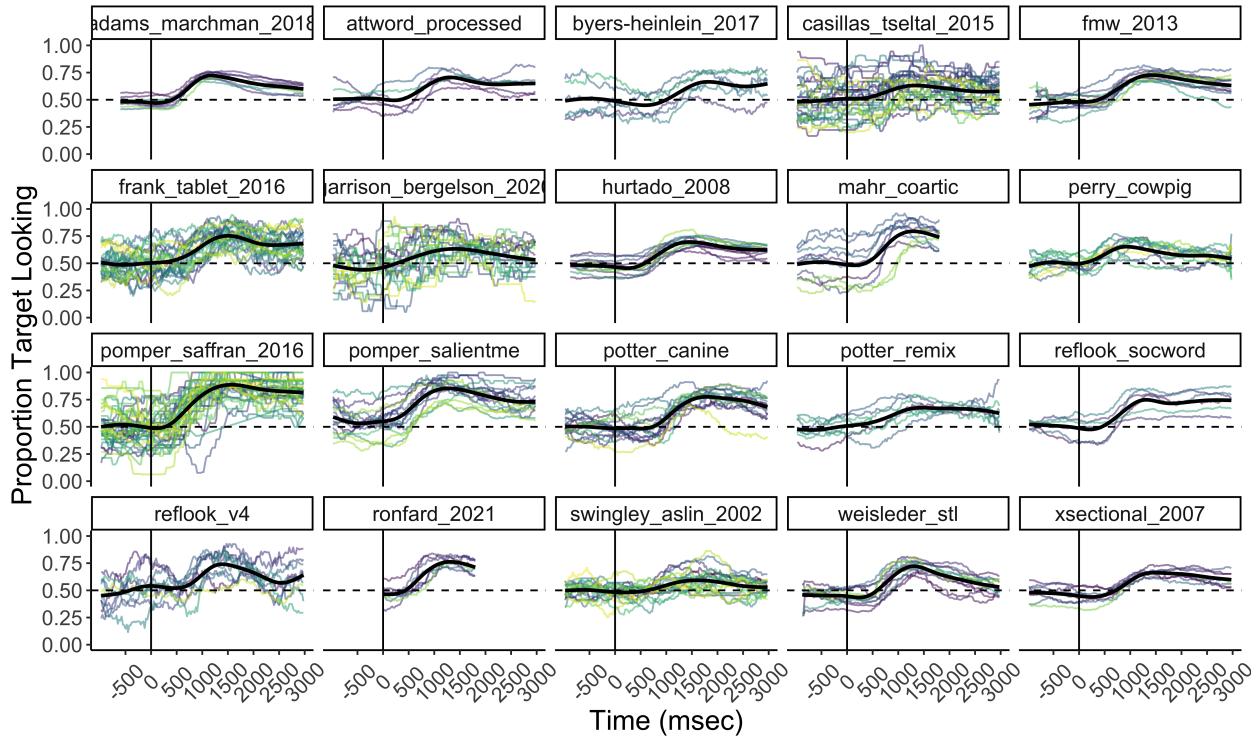
410

### Peekbank: General Descriptives

411 One of the values of the uniform data format we use in Peekbank is the ease of  
 412 providing cross-dataset descriptions that can give an overview of some of the general  
 413 patterns found in our data.

414 A first broad question is about the degree of accuracy in word recognition found across  
 415 studies. In general, participants demonstrated robust, above-chance word recognition in each  
 416 dataset (chance=0.5). Table 2 shows the average proportion of target looking within a  
 417 standard critical window of 367-2000ms after the onset of the label for each dataset  
 418 (Swingley & Aslin, 2002). Proportion target looking was generally higher for familiar words  
 419 ( $M = 0.66$ , 95% CI = [0.65, 0.67],  $n = 1543$ ) than for novel words learned during the  
 420 experiment ( $M = 0.59$ , 95% CI = [0.58, 0.61],  $n = 822$ ).

421 A second question of interest is about the variability across items (i.e., target labels)  
 422 within specific studies. Some studies use a smaller set of items [e.g., 8 nouns; Adams et al.



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

423 (2018)] while others use dozens of different items (e.g., Garrison, Baudet, Breitfeld, Aberman,  
 424 & Bergelson, 2020). Figure 4 gives an overview of the variability in proportion looking to the  
 425 target item for individual words in each dataset. Although all datasets show a gradual rise in  
 426 average proportion target looking over chance performance, the number of unique target  
 427 labels and their associated accuracy vary widely across datasets.

428

## Peekbank in Action

429 We provide two potential use-cases for Peekbank data. In each case, we provide sample  
 430 code to demonstrate the ease of doing simple analyses using the database. Our first example  
 431 shows how we can investigate the findings of a classic study. This type of investigation can  
 432 be a very useful exercise for teaching students about best practices for data analysis (e.g.,

433 Hardwicke et al., 2018) and also provides an easy way to explore looking-while-listening time  
 434 course data in a standardized format. Our second example shows an in-depth exploration of  
 435 developmental changes in the recognition of particular words. Besides its theoretical interest  
 436 (which we will explore more fully in subsequent work), this type of analysis could in principle  
 437 be used for optimizing the stimuli for new experiments, especially as the Peekbank dataset  
 438 grows and gains coverage over a greater number of items.

#### 439 Investigating prior findings: Swingley and Aslin (2002)

440 Swingley and Aslin (2002) investigated the specificity of 14-16 month-olds' word  
 441 representations using the looking-while-listening paradigm, asking whether recognition would  
 442 be slower and less accurate for mispronunciations, e.g. *oppel* (close mispronunciation) or *opel*  
 443 (distant mispronunciation) instead of *apple* (correct pronunciation). In this short vignette,  
 444 we show how easily the data in Peekbank can be used to visualize this result. Our goal here  
 445 is not to provide a precise analytical reproduction of the analyses reported in the original  
 446 paper, but rather to demonstrate the use of the Peekbank framework to analyze datasets of  
 447 this type. In particular, because Peekbank uses a uniform data import standard, it is likely  
 448 that there will be minor numerical discrepancies between analyses on Peekbank data and  
 449 analyses that use another processing pipeline.

```
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")
```

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

454 containing all the information necessary for the analysis.

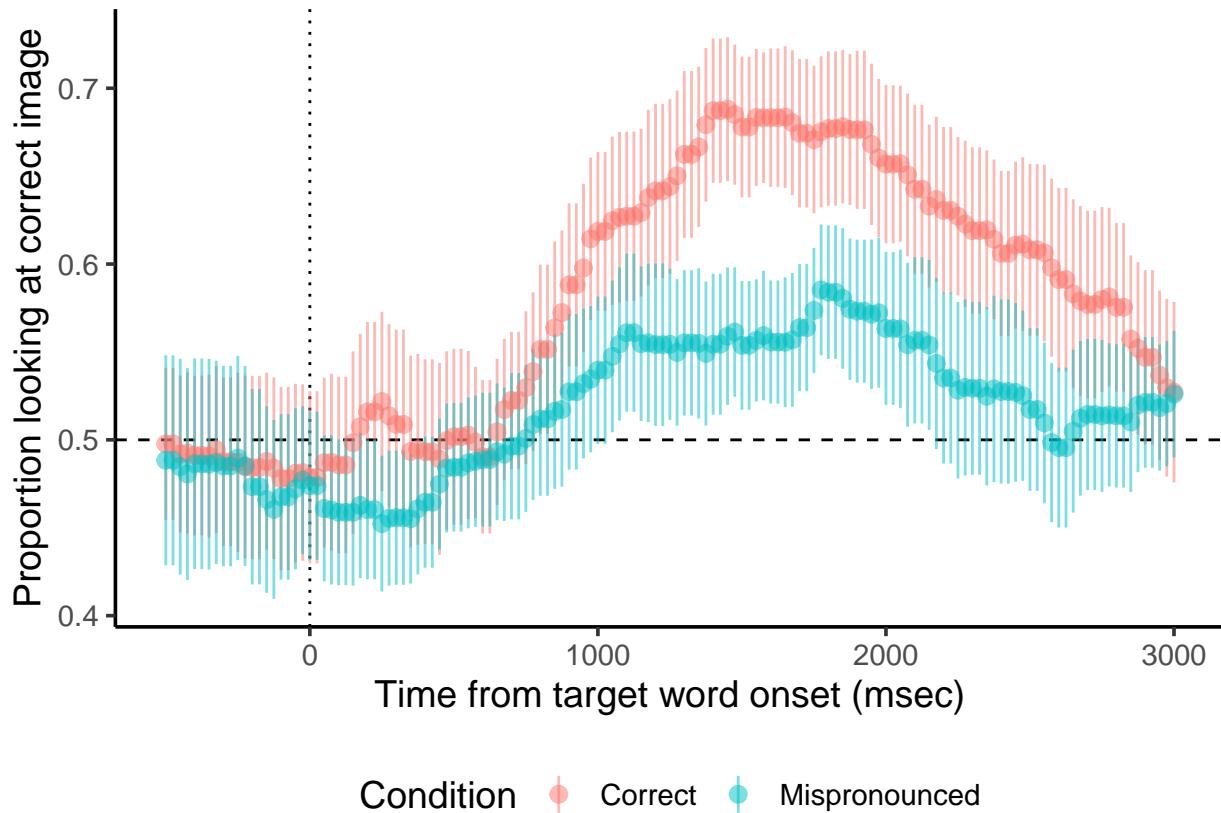
```
swingley_data <- aoi_timepoints |>
  left_join(administrations) |>
  left_join(trials) |>
  left_join(trial_types) |>
  filter(condition != "filler") |>
  mutate(condition = if_else(condition == "cp", "Correct", "Mispronounced"))
```

455 As the code above shows, once the data are joined, condition information for each  
 456 timepoint is present and so we can easily filter out filler trials and set up the conditions for  
 457 further analysis. For simplicity, here we combine both mispronunciation conditions since the  
 458 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()))
```

459 The final step in our analysis is to create a summary dataframe using `dplyr`  
 460 commands. We first group the data by timestep, participant, and condition and compute the  
 461 proportion looking at the correct image. We then summarize again, averaging across  
 462 participants, computing both means and 95% confidence intervals (via the approximation of  
 463 1.96 times the standard error of the mean). The resulting dataframe can be used for  
 464 visualization of the time course of looking.

465 Figure 5 shows the average time course of looking for the two conditions, as produced  
 466 by the code above. Looks after the correctly pronounced noun appeared both faster  
 467 (deviating from chance earlier) and more accurate (showing a higher asymptote). Overall,  
 468 this example demonstrates the ability to produce this visualization in just a few lines of code.



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

#### 469 Item analyses

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

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

478 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)
```

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

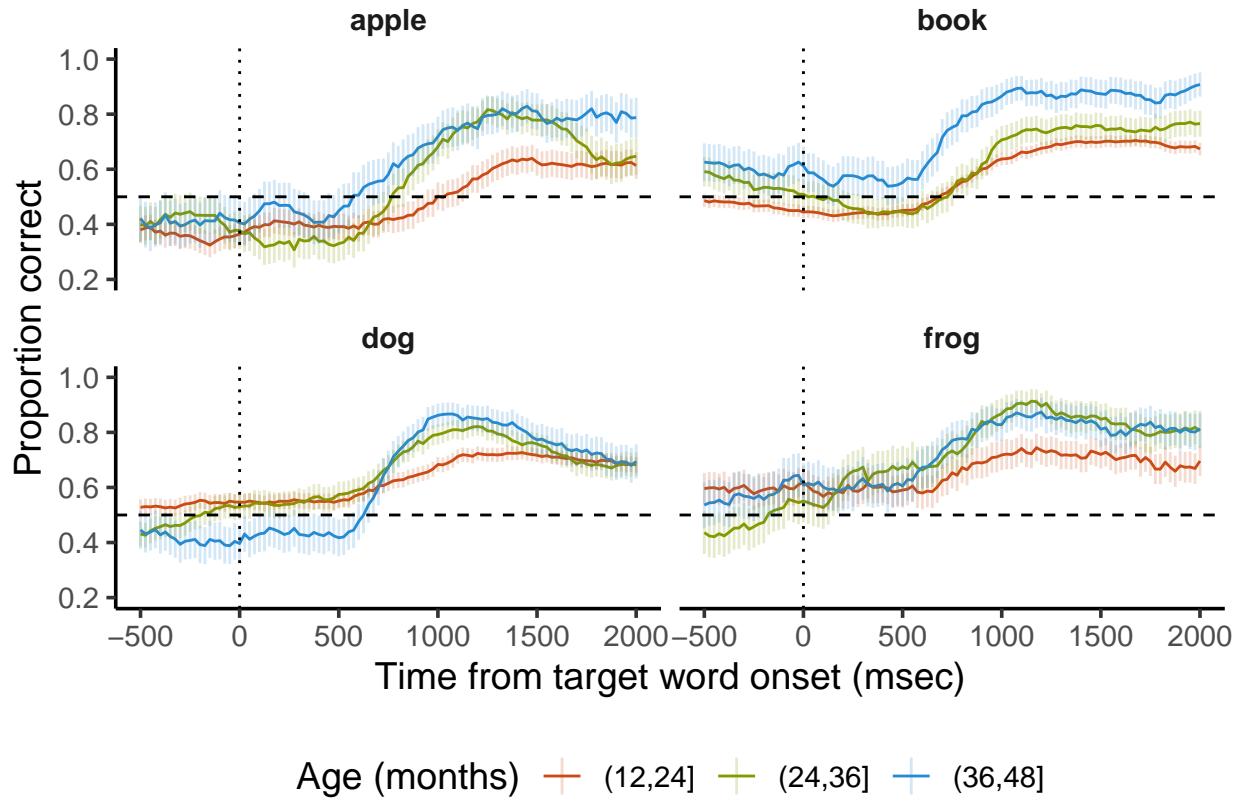
```
target_words <- c("book", "dog", "frog", "apple")

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())

```



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

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

below (with styling commands again removed for clarity). Figure 6 shows the resulting plot, with time courses for each of three (rather coarse) age bins. Although some baseline effects are visible across items, we still see clear and consistent increases in looking to the target, with the increase appearing earlier and in many cases asymptoting at a higher level for older children. On the other hand, this simple averaging approach ignores study-to-study variation

489 (perhaps responsible for the baseline effects we see in the *apple* and *frog* items especially). In  
 490 future work, we hope to introduce model-based analytic methods that use mixed effects  
 491 regression to factor out study-level and individual-level variance in order to recover  
 492 developmental effects more appropriately (see e.g., Zettersten et al., 2021 for a prototype of  
 493 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)
```

494

## Discussion

495 Theoretical progress in understanding child development requires rich datasets, but  
 496 collecting child data is expensive, difficult, and time-intensive. Recent years have seen a  
 497 growing effort to build open source tools and pool research efforts to meet the challenge of  
 498 building a cumulative developmental science (Bergmann et al., 2018; Frank, Braginsky,  
 499 Yurovsky, & Marchman, 2017; Sanchez et al., 2019; The ManyBabies Consortium, 2020).

500 The Peekbank project expands on these efforts by building an infrastructure for aggregating  
 501 eye-tracking data across studies, with a specific focus on the looking-while-listening  
 502 paradigm. This paper presents an overview of the structure of the database, shows how users  
 503 can access the database, and demonstrates how it can be used both to investigate prior  
 504 experiments and to synthesize data across studies.

505 The current database has a number of limitations, particularly in its number and  
 506 diversity of datasets. With 20 datasets currently available in the database, idiosyncrasies of

507 particular designs and condition manipulations still have substantial influence on modeling  
508 results. Expanding the set of distinct datasets will allow us to increase the number of  
509 observations per item across datasets, leading to more robust generalizations across item-level  
510 variability. The current database is also limited by the relatively homogeneous background of  
511 its participants, both with respect to language (almost entirely monolingual native English  
512 speakers) and cultural background (Henrich, Heine, & Norenzayan, 2010; Muthukrishna et  
513 al., 2020). Increasing the diversity of participant backgrounds and languages will expand the  
514 scope of the generalizations we can form about child word recognition.

515 Finally, while the current database is focused on studies of word recognition, the tools  
516 and infrastructure developed in the project can in principle be used to accommodate any  
517 eye-tracking paradigm, opening up new avenues for insights into cognitive development. Gaze  
518 behavior has been at the core of many of the key advances in our understanding of infant  
519 cognition (Bergelson & Swingley, 2012; Fernald, Pinto, Swingley, Weinberg, & McRoberts,  
520 1998; Lew-Williams & Fernald, 2007; Weisleder & Fernald, 2013; Yurovsky & Frank, 2017).  
521 Aggregating large datasets of infant looking behavior in a single, openly-accessible format  
522 promises to bring a fuller picture of infant cognitive development into view.

523 **CRediT author statement**

524 An overview over authorship contributions following the CRediT taxonomy can be  
525 viewed here: [https://docs.google.com/spreadsheets/d/e/2PACX-1vRD-LJD\\_dTAQaAynyBlwXvGpfAVzP-3Pi6JTDoG15m3PYZe0c44Y12U2a\\_hwdmhIstpjyigG2o3na4y/pubhtml](https://docs.google.com/spreadsheets/d/e/2PACX-1vRD-LJD_dTAQaAynyBlwXvGpfAVzP-3Pi6JTDoG15m3PYZe0c44Y12U2a_hwdmhIstpjyigG2o3na4y/pubhtml)  
526

527 **Acknowledgements**

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

530

## References

- 531 Adams, K. A., Marchman, V. A., Loi, E. C., Ashland, M. D., Fernald, A., & Feldman,  
532 H. M. (2018). Caregiver talk and medical risk as predictors of language outcomes  
533 in full term and preterm toddlers. *Child Development*, 89(5), 1674–1690.
- 534 Bergelson, E. (2020). The comprehension boost in early word learning: Older infants  
535 are better learners. *Child Development Perspectives*, 14(3), 142–149.
- 536 Bergelson, E., & Swingley, D. (2012). At 6-9 months, human infants know the  
537 meanings of many common nouns. *PNAS*, 109(9), 3253–3258.
- 538 Bergelson, E., & Swingley, D. (2013). The acquisition of abstract words by young  
539 infants. *Cognition*, 127(3), 391–397.
- 540 Bergmann, C., Tsuji, S., Piccinini, P. E., Lewis, M. L., Braginsky, M., Frank, M. C.,  
541 & Cristia, A. (2018). Promoting replicability in developmental research through  
542 meta-analyses: Insights from language acquisition research. *Child Development*,  
543 89(6), 1996–2009.
- 544 Bleses, D., Makransky, G., Dale, P. S., Højen, A., & Ari, B. A. (2016). Early  
545 productive vocabulary predicts academic achievement 10 years later. *Applied  
546 Psycholinguistics*, 37(6), 1461–1476.
- 547 Byers-Heinlein, K., Bergmann, C., & Savalei, V. (2021). Six solutions for more reliable  
548 infant research. *PsyArXiv*. <https://doi.org/https://doi.org/10.31234/osf.io/ksfvq>
- 549 DeBolt, M. C., Rhemtulla, M., & Oakes, L. M. (2020). Robust data and power in  
550 infant research: A case study of the effect of number of infants and number of  
551 trials in visual preference procedures. *Infancy*, 25(4), 393–419.  
552 <https://doi.org/10.1111/inf.12337>

- 553 Fernald, A., Marchman, V. A., & Weisleder, A. (2013). SES differences in language  
554 processing skill and vocabulary are evident at 18 months. *Developmental Science*,  
555 16(2), 234–248. <https://doi.org/10.1111/desc.12019>
- 556 Fernald, A., Pinto, J. P., Swingley, D., Weinberg, A., & McRoberts, G. W. (1998).  
557 Rapid gains in speed of verbal processing by infants in the 2nd year. *Psychological  
558 Science*, 9(3), 228–231.
- 559 Fernald, A., Zangl, R., Portillo, A. L., & Marchman, V. A. (2008). Looking while  
560 listening: Using eye movements to monitor spoken language comprehension by  
561 infants and young children. In I. A. Sekerina, E. M. Fernandez, & H. Clahsen  
562 (Eds.), *Developmental psycholinguistics: On-line methods in children's language  
563 processing* (pp. 97–135). Amsterdam: John Benjamins.
- 564 Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., ...  
565 Yurovsky, D. (2017). A Collaborative Approach to Infant Research: Promoting  
566 Reproducibility, Best Practices, and Theory-Building. *Infancy*, 22(4), 421–435.  
567 <https://doi.org/10.1111/infa.12182>
- 568 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2017). Wordbank:  
569 An open repository for developmental vocabulary data. *Journal of Child  
570 Language*, 44(3), 677–694.
- 571 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability  
572 and Consistency in Early Language Learning: The Wordbank Project*. Cambridge,  
573 MA: MIT Press.
- 574 Garrison, H., Baudet, G., Breitfeld, E., Aberman, A., & Bergelson, E. (2020).  
575 Familiarity plays a small role in noun comprehension at 12–18 months. *Infancy*,  
576 25(4), 458–477.

577 Gautheron, L., Rochat, N., & Cristia, A. (under review). Managing, storing, and  
578 sharing long-form recordings and their annotations. Retrieved from  
579 <https://doi.org/10.31234/osf.io/w8trm>

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

583 Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P.,  
584 ... Poldrack, R. A. (2016). The brain imaging data structure, a format for  
585 organizing and describing outputs of neuroimaging experiments. *Scientific Data*,  
586 3(1), 160044. <https://doi.org/10.1038/sdata.2016.44>

587 Hardwicke, T. E., Mathur, M. B., MacDonald, K., Nilsonne, G., Banks, G. C.,  
588 Kidwell, M. C., ... Frank, M. C. (2018). Data availability, reusability, and  
589 analytic reproducibility: Evaluating the impact of a mandatory open data policy  
590 at the journal Cognition. *Royal Society Open Science*, 5(8).  
591 <https://doi.org/10.1098/rsos.180448>

592 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world?  
593 *The Behavioral and Brain Sciences*, 33(2-3), 61–83.  
594 <https://doi.org/10.1017/S0140525X0999152X>

595 Hirsh-Pasek, K., Cauley, K. M., Golinkoff, R. M., & Gordon, L. (1987). The eyes  
596 have it: Lexical and syntactic comprehension in a new paradigm. *Journal of Child  
597 Language*, 14(1), 23–45.

598 Hurtado, N., Marchman, V. A., & Fernald, A. (2007). Spoken word recognition by  
599 Latino children learning Spanish as their first language. *Journal of Child  
600 Language*, 34(2), 227–249. <https://doi.org/10.1017/S0305000906007896>

- 601 Hurtado, N., Marchman, V. A., & Fernald, A. (2008). Does input influence uptake?  
602 Links between maternal talk, processing speed and vocabulary size in  
603 Spanish-learning children. *Developmental Science*, 11(6), 31–39.  
604 <https://doi.org/10.1111/j.1467-7687.2008.00768.x>
- 605 Lewis, M., Braginsky, M., Tsuji, S., Bergmann, C., Piccinini, P. E., Cristia, A., &  
606 Frank, M. C. (2016). *A Quantitative Synthesis of Early Language Acquisition  
Using Meta-Analysis*. <https://doi.org/10.31234/osf.io/htsjm>
- 608 Lew-Williams, C., & Fernald, A. (2007). Young children learning Spanish make rapid  
609 use of grammatical gender in spoken word recognition. *Psychological Science*,  
610 18(3), 193–198.
- 611 Mahr, T., McMillan, B. T. M., Saffran, J. R., Ellis Weismer, S., & Edwards, J. (2015).  
612 Anticipatory coarticulation facilitates word recognition in toddlers. *Cognition*,  
613 142, 345–350. <https://doi.org/10.1016/j.cognition.2015.05.009>
- 614 Marchman, V. A., Loi, E. C., Adams, K. A., Ashland, M., Fernald, A., & Feldman, H.  
615 M. (2018). Speed of language comprehension at 18 months old predicts  
616 school-relevant outcomes at 54 months old in children born preterm. *Journal of  
617 Dev. & Behav. Pediatrics*, 39(3), 246–253.
- 618 Muthukrishna, M., Bell, A. V., Henrich, J., Curtin, C. M., Gedranovich, A.,  
619 McInerney, J., & Thue, B. (2020). Beyond Western, Educated, Industrial, Rich,  
620 and Democratic (WEIRD) Psychology: Measuring and Mapping Scales of  
621 Cultural and Psychological Distance. *Psychological Science*, 31(6), 678–701.
- 622 Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A.,  
623 ... Vazire, S. (2021). Replicability, Robustness, and Reproducibility in  
624 Psychological Science. *PsyArXiv*.

- 625 https://doi.org/https://doi.org/10.31234/osf.io/ksfvq
- 626 Peter, M. S., Durrant, S., Jessop, A., Bidgood, A., Pine, J. M., & Rowland, C. F.  
627 (2019). Does speed of processing or vocabulary size predict later language growth  
628 in toddlers? *Cognitive Psychology*, 115, 101238.
- 629 Potter, C., & Lew-Williams, C. (unpublished). Behold the canine!: How does  
630 toddlers' knowledge of typical frames and familiar words interact to influence their  
631 sentence processing?
- 632 R Core Team. (2020). *R: A language and environment for statistical computing*.  
633 Vienna, Austria: R Foundation for Statistical Computing. Retrieved from  
634 <https://www.R-project.org/>
- 635 Ronfard, S., Wei, R., & Rowe, M. L. (2021). Exploring the linguistic, cognitive, and  
636 social skills underlying lexical processing efficiency as measured by the  
637 looking-while-listening paradigm. *Journal of Child Language*, 1–24.  
638 <https://doi.org/10.1017/S0305000921000106>
- 639 Sanchez, A., Meylan, S. C., Braginsky, M., MacDonald, K. E., Yurovsky, D., & Frank,  
640 M. C. (2019). childe-db: A flexible and reproducible interface to the child  
641 language data exchange system. *Behavior Research Methods*, 51(4), 1928–1941.  
642 <https://doi.org/10.3758/s13428-018-1176-7>
- 643 Swingley, D., & Aslin, R. N. (2002). Lexical neighborhoods and the word-form  
644 representations of 14-month-olds. *Psychological Science*, 13(5), 480–484.  
645 <https://doi.org/10.1111/1467-9280.00485>
- 646 The ManyBabies Consortium. (2020). Quantifying sources of variability in infancy  
647 research using the infant-directed speech preference. *Advances in Methods and*  
648 *Practices in Psychological Science*, 3(1), 24–52.

- 649 Wass, S. V., Smith, T. J., & Johnson, M. H. (2013). Parsing eye-tracking data of  
650 variable quality to provide accurate fixation duration estimates in infants and  
651 adults. *Behavior Research Methods*, 45(1), 229–250.  
652 <https://doi.org/10.3758/s13428-012-0245-6>
- 653 Weisleder, A., & Fernald, A. (2013). Talking to Children Matters: Early Language  
654 Experience Strengthens Processing and Builds Vocabulary. *Psychological Science*,  
655 24(11), 2143–2152. <https://doi.org/10.1177/0956797613488145>
- 656 Yurovsky, D., & Frank, M. C. (2017). Beyond naïve cue combination: salience and  
657 social cues in early word learning. *Dev Sci*, 20(2).
- 658 Yurovsky, D., Wade, A., Kraus, A. M., Gengoux, G. W., Hardan, A. Y., & Frank, M.  
659 C. (under review). Developmental changes in the speed of social attention in early  
660 word learning.
- 661 Zettersten, M., Bergey, C., Bhatt, N., Boyce, V., Braginsky, M., Carstensen, A., ...  
662 others. (2021). Peekbank: Exploring children's word recognition through an open,  
663 large-scale repository for developmental eye-tracking data.