

1 Developmental changes in the precision of visual concept knowledge

2 Bria Long¹ & Ernst-August Doelle^{1,2}

3 ¹ University of California San Diego

4 ² Stanford University

5 Author Note

6 Add complete departmental affiliations for each author here. Each new line herein
7 must be indented, like this line.

8 Enter author note here.

9 The authors made the following contributions. Bria Long: Conceptualization,
10 Writing - Original Draft Preparation, Writing - Review & Editing; Ernst-August Doelle:
11 Writing - Review & Editing, Supervision.

12 Correspondence concerning this article should be addressed to Bria Long, 9500
13 Gillman Drive, La Jolla, CA 92093. E-mail: brlong@ucsd.com

Abstract

How precise is children’s visual concept knowledge, and how does this change across development? We created a gamified picture-matching task where children heard a word (e.g., “swordfish”) and had to choose the picture “that goes with the word.” We collected data from large sample of children on this task, and we modeled changes in the proportion of children who chose a given image for a certain word over development. We found gradual developmental changes in children’s ability to identify the correct category. Error analysis revealed that children were more likely to choose higher-similarity distractors as they grew older; children’s error patterns were increasingly correlated with CLIP target-distractor similarity. These analyses suggest a transition from coarse to finer-grained visual representations over early and middle childhood. Broadly, these findings demonstrate the utility of combining gamified experiments and similarity estimates from computational models to probe the content of children’s evolving visual representations.

Keywords: visual concepts, receptive vocabulary, large language models, object recognition

Word count: X

Developmental changes in the precision of visual concept knowledge

```
## Warning: package 'kableExtra' was built under R version 4.3.2
```

Introduction

When a child hears a word — like a “whale” — this activates a mental representation of its referent in the real-world. But what is this representation actually like? Depending on how old a child is—and how much they have learned about whales—they might imagine a canonical exemplar of a blue whale, a specific whale from a picture book, or perhaps they just know vaguely that a whale is an animal that lives in the ocean. How precise are the visual representations that underlie children’s understandings of words across childhood?

As infants begin to communicate with their caregivers, they experience an astonishing rate of vocabulary growth (Bloom, 2000; Braginsky et al., 2022). Infants as young as 6-months of age appear to absorb some shape information from label-object co-occurrences in everyday experience (Vong et. al 2024; Bergelson et al., 2009), and 14-18 month-olds extend newly learned words to atypical exemplars (Weaver et al., 2024, Child Dev). By around their second birthday, children extend words to stylized, 3D exemplars (Smith, 2003) as they learn that shape is a valuable cue to basic-level categories (Rosch et al., 1976). Thus, at least for within-category exemplars, very young children exhibit relatively sophisticated generalization abilities for common visual concepts, in line with a broad-to-narrow view of category development (Waxman & Gelman, 2009), where infants construe words as initially referring to many items and subsequently refine their representations across development. From this perspective, children’s visual representations may change relatively little across childhood; instead, children may gradually acquire new visual concepts and instead change in how they represent the relationships between visual concepts: for example, children may learn that whales are mammals, and then appropriately group them with other land mammals vs. with fish when asked to make

55 taxonomic classifications. Accordingly, empirical work on children’s developing ability to
56 recognize objects (Azyenberg & Behrman, 2024) has also focused on the first few years of
57 childhood as the most critical period in which object recognition abilities develops.

58 Contrary to this simplified account, here we posit that children’s visual concepts
59 change throughout childhood, with an extended developmental trajectory that continues in
60 parallel with later vocabulary learning. Of course, children’s vocabulary knowledge—often
61 assessed via paper-and-pencil, closed, expensive traditional assessments—grows and
62 expands across childhood (CITE), but there has been relatively little consideration of the
63 visual representations that support children’s performance on picture vocabulary tasks. As
64 children enter schooling environments and begin to learn why animals and objects are
65 classified the way they are, this semantic learning is likely to influence the visual features
66 that are prioritized in children’s visual concepts. Some work on children’s production and
67 recognition of drawings of common objects hints at this kind of protracted developmental
68 timeline (Long et al., 2024): in a large observational study, children became increasingly
69 able to both depict and recognize line drawings of common object categories. However, no
70 work has directly tested children’s visual recognition behaviors for a wide variety of visual
71 concepts. We suspect that this is in part because of the difficulty of obtaining data from
72 large samples of children on a consistent set of items with variability over a large
73 developmental age range.

74 To overcome these methodological barriers, we created a gamified picture-matching
75 task where children heard a word (e.g., “swordfish”) and had to choose the picture “that
76 goes with the word”. Critically, we chose distractor items with high, medium, and low
77 concept similarity to each target word; distractors were paired via cosine similarity of the
78 target and distractor words in a large multimodal language model (CLIP, Radford et al.,
79 2021). This task was then deployed in online, preschool, and school contexts to 3599
80 children aged 3-15 years and 211 adults years of age. Using this large dataset, we find
81 gradual changes in how children represent visual concepts across childhood, with older

children becoming both more accurate at identifying the correct referents throughout this extended age range; however, we also found that even young children were more likely to choose the related vs. unrelated distractors, highlighting a gradual change from coarse, representations that encompass both the target and related distractors to fine-grained, specific representations that the visual information that words refer to. We then use both unimodal and multimodal embeddings from this same model to examine how visual, linguistic, and multimodal similarity explain children’s error patterns across development.

Methods

Procedure

Children were invited to participate in a picture matching game; a cover story accompanied the game where children were asked to help teach aliens on another planet about some of the words on our planet; children were able to pick a particular alien to “accompany them on their journey.” Before the stimuli appeared, children heard a target word (e.g., “apple”) and then were asked to “choose the picture that goes with the word”. The four images appeared in randomized locations on the screen, and one of the images always corresponded to the target word. On practice trials, the distractor images were all very dissimilar to the target concept, and the target word was relatively easy. The tablet played a chime sound if they chose correctly, and a slightly unpleasant sound if they responded incorrectly. Each child viewed a random subset of the item bank, and the items they viewed were displayed in a random order. Children were allowed to stop the game if they wanted to. While different versions of the game included varying amounts of trials or items, as these games are part of a larger project to develop an open-sourced measure of children’s vocabulary knowledge as an alternative to the PPVT; however, here we analyze children’s responses to items that were generated using the THINGS+ dataset with distractors of varying difficulty.

Participants

To obtain a large sample, we collected data from children in several different testing contexts. We collected data from children in an in-person preschools ($N = 65$, 3-5 year-olds), from the Children Helping Science Platform, ($N=243$, 3-7 year-olds), 6 elementary schools, and 9 charter schools across multiple states ($N=3332$, 5-14 year-olds) and adults online ($N=211$ adults, recruited via Prolific; half of the adults spoke English as a second language). Most participants responded via a touch-screen tablet, except those recruited online: however, children’s parents responded via clicking on the image on Children Helping Science, and adults responded via clicking on the images.

After pre-processing, we included for a total of 3786 participants from preschools, schools, and online testing contexts around the country (range 84 to 654), who completed, on average, 25.02 4AFC trials that were sampled randomly from the stimuli set (max=86; different maximum numbers of trials were included in different testing contexts). We tested an additional 84 participants who scored near chance on 4AFC trials (chance=25%, threshold=30%) and were school-aged (>6 years of age) and who we excluded from analyses; these participants completed an average of 17.72 trials.

Stimuli selection

We capitalized on publicly available existing image and audio databases to generate stimuli. Visual concepts were taken from the THINGS+ dataset (Stoinski et al., 2023), after filtering out non-child safe images (e.g., weapons, cigarettes) and images with low nameability, as per the released norming data. We used the copy-right free, high-quality image released for each visual concept. We then subset to visual concepts that had available audio recordings in the MALD database as well as age-of-acquisition (AoA) ratings from a previous existing dataset (Kuperman, 2012).

Using this subset, we sampled distractors with high, medium, and low similarity to

the target word as operationalized via embedding similarity of the words in the language encode of a multimodal large language model (CLIP, Contrastive Language-Image Pre-training, Radford et al., 2021). High-, medium, and low similarity values were determined relative to the distribution of possible target-distractor pairing values for each word in the THINGS+ dataset. In our final set, we had 108 items with a range of different estimated age-of-acquisitions (e.g., hedgehog, mandolin, mulch, swordfish, waterwheel, bobsled) with all unique targets and distractors; in addition, we constrained the sampling such that target-distractor pairs had estimated age of acquisition within 3 years of each other. All stimuli and their meta-data are available on the public repository for this project.

Model features

We obtained all model features using the OpenAI available implementation of CLIP available at <https://github.com/openai/CLIP>. For language similarity, we computed the cosine similarity of the embeddings of the target to each distractor word on each trial (e.g., rose – tulip, rose – glove, rose – hubcap). For visual similarity, we repeated this procedure but by obtaining image similarity vectors in the vision transformer. For multimodal similarity, we computed the cosine similarity of the embedding of the target word in the language model to the embeddings for each of the distractor images; this is possible because the embedding spaces for the vision and language transformers in the CLIP model are aligned and have the same number of dimensions.

Results

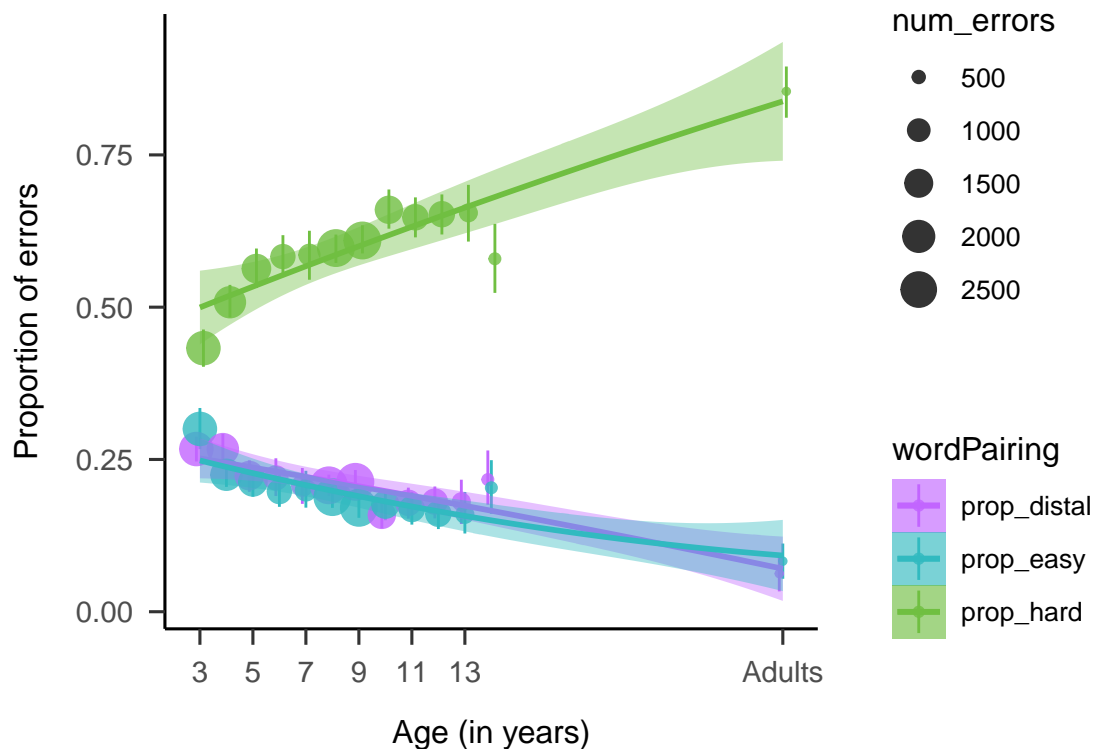
A protracted developmental trajectory

We found a gradual increase in children’s ability to correctly identify the target word across our entire age range, extending into early adolescence; Figure ?? shows the proportion of time that children identified the target word, highlighting a protracted

developmental trajectory. We found this developmental trend for both relatively “easy” words, with an average estimated age-of-acquisition (AoA) of 4.81 years ($SD = 0.87$), more difficult words (Average AoA = 6.95 years, $SD = 0.65$ years), and challenging words (Average AoA = 9.60 years, $SD = 1.21$)).

At an item level, the words that showed the greatest change across age included some animals (e.g., “swordfish”) as well as inanimate objects (“prism”, “antenna”, “gutter”, “sandbag”, “turbine”) but also parts of larger buildings (“scaffolding”, “gutter”). However, some of our developmental trends likely also stem from differences in executive control: for some words that had very semantically similar distractors but were relatively easy (e.g., cheese and butter), we still saw steep developmental changes, highlighting that this “simple” picture vocabulary matching tasks still assess many different cognitive abilities beyond the fidelity of children’s visual representations.

Changes in the precision of visual concepts



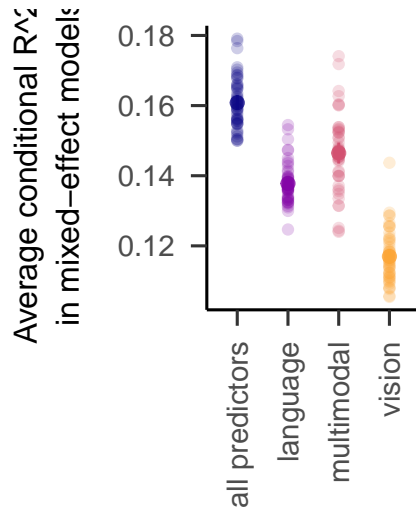
Next, we thus aimed to understand whether we were observing real changes in the precision of children’s visual concepts, or perhaps simply changes in task performance as children become more accurate at ignoring relevant distractors. To examine this, we focused on changes in error patterns across age, shown in Figure ?? . If children’s visual concepts are changing—proceeding from a coarse representations that is overly broad, or starting from no representation at all—then we should observe changes in how children choose distractors when they make errors. Specifically, we would expect younger children to be more likely to choose distractors of all types, whereas we would expect older children to almost exclusively choose related distractors. Consistent with this hypothesis, we found that children increasingly choose related distractors throughout development, with adults being still more likely to choose the related distractors relative to the oldest children (14-year-olds) in our sample.

We confirmed this result via a linear regression, modeling the proportion of errors where children chose the related distractor as our dependent variable as a function of

children’s age (in years), finding a fixed effect of age (see Table 1): older children were more likely to choose related distractors relative to unrelated distractors. For robustness, we also modeled these effects at the item level in a linear mixed-effect model, with random intercepts for each item, finding the same pattern of effects. Children become more likely to choose related distractors across development, suggesting a progression where children gradually build detailed knowledge about the visual referents of many challenging words.

Modeling changes in children’s error patterns

Warning: 'r.squaredGLMM' now calculates a revised statistic. See the help page.



In a final analyses, we aimed to understand the source of these changes in children’s error patterns by leveraging high-dimensional embeddings of our linguistic and visual stimuli in a large, multimodal language model (CLIP, Radford et al., 2021). We chose a set of stimuli where visual similarity was colinear with semantic similarity to a large degree, as it often is in the real-world for most visual concepts. Thus, our stimuli were not necessarily designed to pull apart the contributions of changes in semantic vs. visual similarity. Nonetheless, our stimuli were generated by using similarity in a linguistic embedding space, and so some stimuli on certain trials were nonetheless related to the target concept semantically but not necessarily visually (e.g, gardening gloves were a distractor for

“rose”). We thus sought to understand the degree to which children’s error patterns in this task reflected changes in how they processed the visual similarity of the targets and distractors, their semantic similarity, or—perhaps most likely—some combination.

To do so, our approach used a series of cross-validated linear mixed effect models, where we examined the degree to which visual, linguistic, multimodal, and a combination of similarity metrics derived by large language model embeddings could predict the errors that children made. Overall, we found that...

Discussion

How precise is children’s visual concept knowledge, and how does this change across development?

Overall, these analyses suggest a transition from coarse to finer-grained visual representations over early and middle childhood.

Children’s visual concept knowledge gradually becomes more refined as children learn what distinguishes similar visual concepts from one another. Broadly, these findings demonstrate the utility of combining gamified experiments and similarity estimates from computational models to probe the content of children’s evolving visual representations.

Implications: Supports Ecological enrichment accounts:

On another viewpoint, there is also substantial enrichment and change in children’s visual representations of everyday visual concepts. Broader view on the learning environment (e.g., Bruner), and children as quite active participants in their learning environment, Longer view on the timeline for learning, which in turn changes how we think about the relevant learning environment— which changes substantially as they grow and learn both from their peers throughout early childhood and in structured educational contexts. For example, they may have grossly misrepresented the sizes of certain objects (e.g., whales are XX bigger than dolphins) and certain visual features may become more

or less salient as they understand their functional roles (XX) or semantic relevance of the category. On this account, even school-aged children’s visual representations may undergo substantial change as they learn more about the world around them, even as their vocabulary growth tapers. Goes beyond acquisition account to suggest that their representations change beyond what has been measure in classic recognition tasks with young children]

Connection to adult expertise We suspect that visual concept learning extends into adulthood, and that many adults have coarse visual representations for many different words. Consider that while we experience the referents of some visual concepts relatively frequently—e.g., trees, computers, cups, cars—other words refer to visual concepts that different individuals may have varying amounts of interest in and frequency in interacting with—like telescopes, or antelopes. Visual concept learning is likely influenced by both children and adults’ occupation and pre-occupations. And indeed decades of work has established that birding experts, car aficionados, and graphic artists have both qualitatively and quantitatively different kinds of visual representations for the visual concepts that they engage with (CITE, CITE):

References

Table 1
Fixed Effects from Linear Mixed Effects Model

effect	Predictor	*b*	*SE*	*t*	df
fixed	Intercept	0.62	0.01	122.41	3407.593
fixed	Age (scaled)	0.06	0.01	11.26	3419.562
fixed	Number of trials (scaled)	0.00	0.00	-0.42	3404.545

Note.
Analysis conducted using a linear mixed effects model. The model included random intercepts for partic

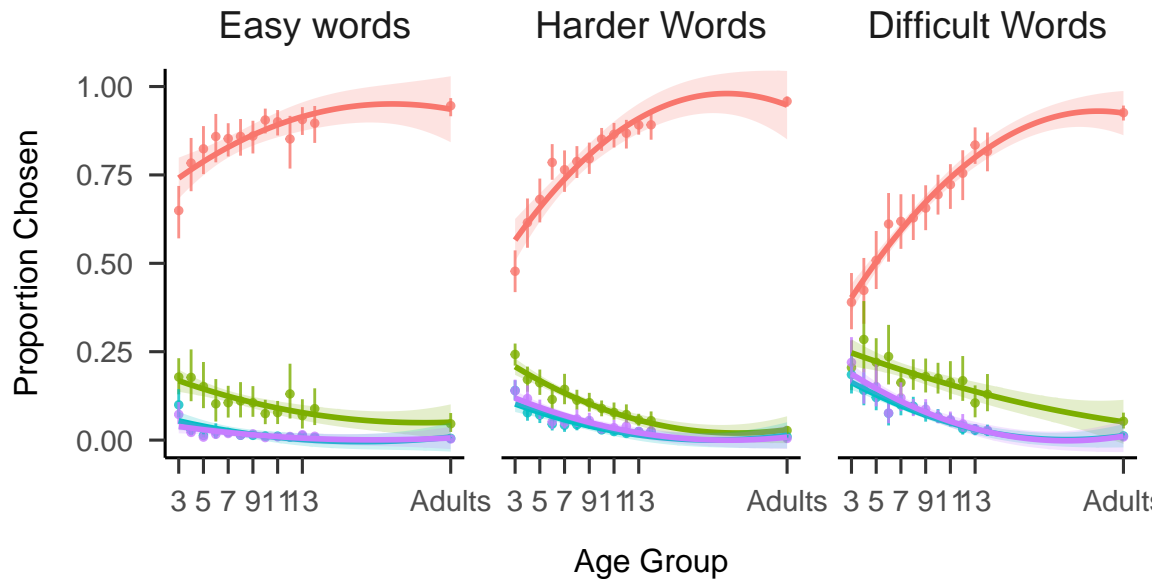


Figure 1. Visual vocabulary task performance as a function of the age of the child completing the task, plotted separately for relatively easy, harder, or difficult words; words are binned into terciles based on the estimated AoA from Kuperman et al., 2012. Lines refer to the proportion of words that children chose the target (red), high-similarity (green), medium similarity (turquoise), or low similarity (purple) distractor at each age; error bars represent bootstrapped confidence intervals.