

Parallel developmental changes in children's production and recognition of  
line drawings of visual concepts

Bria Long<sup>1</sup>, Judith E. Fan<sup>1,2</sup>, Holly Huey<sup>2</sup>, Zixian Chai<sup>1</sup>, Michael C. Frank<sup>1</sup>

<sup>1</sup>Department of Psychology, Stanford University

<sup>2</sup>Department of Psychology, University of California San Diego

## Abstract

Childhood is marked by the rapid accumulation of knowledge and the prolific production of drawings. We conducted a systematic study of how children create and recognize line drawings of visual concepts. We recruited 2-10-year-olds to draw 48 categories via a kiosk, resulting in >37K drawings. We analyzed changes in the category-diagnostic information in these drawings using vision algorithms and annotations of object parts. We found developmental gains in children's inclusion of category-diagnostic information that was not reducible to variation in visuomotor control or effort. Moreover, unrecognizable drawings contained information about the animacy and size of the category children tried to draw. Using "guessing games" at the same kiosk, we found that children improved across childhood at recognizing each other's line drawings. This work leverages vision algorithms to characterize developmental changes in a large dataset of children's drawings and suggests that changes in children's drawings reflect refinements in internal representations.

# Parallel developmental changes in children's production and recognition of line drawings of visual concepts

## Introduction

2 What makes a drawing of a *rabbit* look like a *rabbit* and not a *dog*? As adults, our  
3 *visual concepts* – our sense of what particular objects look like – are seamlessly  
4 integrated into our visual experience. With a single glance, incoming patterns of light  
5 make contact with our visual concepts, supporting the rapid categorization of a wide  
6 variety of inputs, from real-life exemplars to sparse line drawings and stylized  
7 animations (Biederman & Ju, 1988; Gibson, 1971; Hertzmann, 2020; Sayim, 2011). We  
8 can also access our visual concepts in the absence of perceptual input – going beyond  
9 what we have previously experienced to imagine new visual entities and create external  
10 representations of them (Clottes, 2008; Finke & Slayton, 1988; Gregory, 1973).

11 Yet while these feats of perceiving and creating can feel effortless, the  
12 representations that support them are acquired and refined gradually as children learn  
13 about the visual world (Rosch, 1978). Children begin building visual concepts in  
14 earnest during the second year of life as they learn which labels refer to both depictions  
15 and real-life exemplars of categories (DeLoache, Pierroutsakos, & Uttal, 2003). And by  
16 their second birthday, children can learn category labels for novel objects after exposure  
17 to just one or a few exemplars (Carey & Bartlett, 1978; Pereira & Smith, 2009; Soja,  
18 Carey, & Spelke, 1991) and succeed even for sparse 3D representations of these objects  
19 devoid of color and texture-based cues (Pereira & Smith, 2009).

20 But children take many years to learn how to appropriately generalize and  
21 discriminate between visual concepts. For example, children gradually improve in their  
22 ability to accurately group together categories based on taxonomy versus salient  
23 perceptual features (e.g., grouping a *snake* with a *lizard* vs. a *hose*) (Fisher, Godwin, &  
24 Matlen, 2015; Tversky, 1985). Further, children's visual recognition abilities also have a  
25 protracted developmental trajectory throughout middle childhood (Bova et al., 2007;  
26 Juttner, Wakui, Petters, & Davidoff, 2016; Nishimura, Scherf, & Behrmann, 2009) as  
27 children become steadily better at discriminating between similar exemplars of scenes.

28 objects, bodies, and faces between 5-10 years of age (Weigelt et al., 2014), and  
29 increasingly skilled at recognizing objects across unusual poses or 3D rotations, reaching  
30 adult-level performance only in adolescence (Bova et al., 2007; Dekker, Mareschal,  
31 Sereno, & Johnson, 2011; Nishimura, Scherf, Zachariou, Tarr, & Behrmann, 2015). In  
32 turn, changes in children's recognition abilities are related to changes in how the visual  
33 cortex encodes different objects and scenes (Balas & Saville, 2020; Cohen et al., 2019;  
34 Dekker et al., 2011; Gomez, Natu, Jeska, Barnett, & Grill-Spector, 2018; Kersey, Clark,  
35 Lussier, Mahon, & Cantlon, 2015; Nishimura et al., 2015); for example, children's  
36 ability to discriminate similar faces is correlated with the sensitivity of face-selective  
37 regions to these particular faces (Natu et al., 2016). These changes in children's ability  
38 to discriminate exemplars may be driven by children's increasing attention to the  
39 relationships between object parts and their overall configuration (Juttner, Muller, &  
40 Rentschler, 2006; Juttner et al., 2016; Mash, 2006). Together, these findings suggest  
41 that visual concepts are refined throughout childhood as children's perceptual abilities  
42 mature and children learn how to discriminate between similar categories.

43 Psychologists have typically probed children's visual concepts by asking children  
44 to make discrete choices between small samples of stimuli that vary along dimensions  
45 chosen by an experimenter. While valuable for testing specific hypotheses, this strategy  
46 is also characterized by severe limits on the amount of information that can be acquired  
47 on any given experimental trial. By contrast, generative tasks such as drawing  
48 production can overcome these limits by enabling the collection of more information  
49 about the contents of children's visual concepts on every trial. Such tasks are feasible to  
50 administer in experimental settings given that almost all children prolifically produce  
51 drawings of visual concepts from an early age (Piaget, 1929). And there is substantial  
52 precedent for examination of children's drawings to probe their knowledge about the  
53 visual world (Fury, Carlson, & Sroufe, 1997; Karmiloff-Smith, 1990; Kellogg, 1969;  
54 Piaget, 1929). Freehand drawing production tasks thus provide a valuable tool for  
55 characterizing developmental changes in visual concepts. Here we create a large digital  
56 dataset of children's drawings and leverage innovations in machine learning to

57 characterize how changes in children's drawings are related to their growing  
58 understanding of various visual concepts.

59 **Drawings as a window into visual representations**

60 Our work builds on a long literature that has argued that children's drawings of  
61 objects reflect not only what they can directly observe, but what they *know* about these  
62 objects (see "intellectual realism" in Freeman & Janikoun, 1972; Luquet, 1927). For  
63 example, even when drawing from observation, children tend to include features that  
64 are not visible from their vantage point but are nevertheless diagnostic of category  
65 membership (e.g., an occluded handle on a *mug*) (Barrett & Light, 1976; Bremner &  
66 Moore, 1984). Further, direct visual or haptic experience with novel objects tends to  
67 change what information children draw (Bremner & Moore, 1984). These initial studies  
68 have focused on a small number of visual concepts – especially the human figure  
69 (Goodenough, 1963) – finding that younger children (4-5 years) tend to include fewer  
70 category-diagnostic cues in their drawings, such as those distinguishing an "adult" from  
71 a "child," than somewhat older children (6 years), who tend to enrich their drawings  
72 with more diagnostic part information (Cox & Ralph, 1996; Sitton & Light, 1992).  
73 However, the generality of the conclusions based on this work has been unclear given  
74 the narrow range of concepts tested and the lack of generic methods for measuring  
75 diagnostic information in drawings. Further, little work has systematically related  
76 children's ability to include diagnostic visual information in drawings to their emerging  
77 abilities to control and plan their motor movements – which certainly influence how and  
78 what children draw (Freeman, 1987; Rehrig & Stromswold, 2018).

79 Yet research in adults does suggest that what we draw is tightly linked to what we  
80 know about objects and how we perceive them. For example, patients with semantic  
81 dementia tend to produce drawings without distinctive visual features (Bozeat et al.,  
82 2003) or include erroneous features (e.g., a duck with four legs). One recent study found  
83 that adults can produce detailed drawings of scenes after only viewing them for a few  
84 seconds, interleaved among other scenes (Bainbridge, Hall, & Baker, 2019). Another

85 study found that recognizing an object and producing a drawing of an object recruit a  
86 shared neural representation in early visual cortex (Fan et al., 2020). Further, practice  
87 producing drawings of objects can impact perceptual judgments about them. In one  
88 study, adult participants who repeatedly drew similar objects (i.e. *beds* vs *chairs*) were  
89 better able to distinguish them in a categorization task (Fan, Yamins, & Turk-Browne,  
90 2018). Drawing expertise is also associated with enhanced visual encoding of object  
91 parts and their relationships (Chamberlain et al., 2019; Perdreau & Cavanagh, 2013a,  
92 2013b, 2014), but not differences in low-level visual processing (Chamberlain et al.,  
93 2019; Chamberlain, Kozbelt, Drake, & Wagemans, 2021; Kozbelt, 2001; Perdreau &  
94 Cavanagh, 2013b) or shape tracing skills (Tchalenko, 2009).

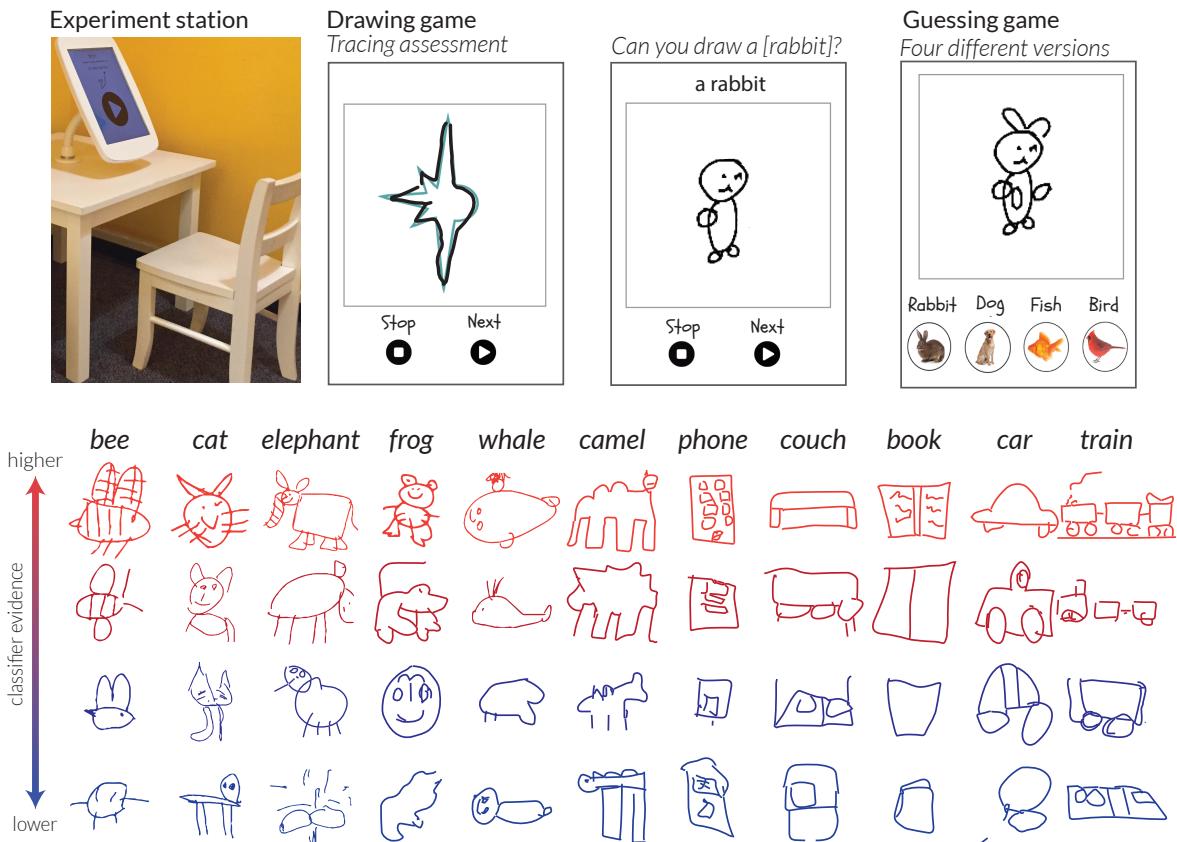
## 95 **The current study**

96 Building on these traditions, in the current paper we characterize developmental  
97 changes in how children produce and recognize line drawings as an additional lens into  
98 children's growing understanding of these visual concepts. We anticipated that  
99 children's ability to produce and recognize line drawings would continue to develop  
100 beyond the preschool and elementary school years (Dekker et al., 2011; Gomez et al.,  
101 2018; Weigelt et al., 2014) and that some – but not all – age-related variation in  
102 drawing ability would be due to improvements in planning and motor control (Freeman,  
103 1987; Rehrig & Stromswold, 2018). In particular, as children learn the visual  
104 information most diagnostic of a visual concept (Rosch, 1978), this visual knowledge  
105 may manifest in both: (1) an enhanced ability to produce line drawings that contain  
106 category-diagnostic information and (2) a greater sensitivity to this same visual  
107 information when recognizing line drawings made by other children.

108 We thus first collected digital drawings of 48 different visual concepts from a large  
109 sample of children spanning a wide age range (2-10 years), resulting in a corpus  
110 containing >37K drawings. To quantify developmental changes in these drawings at  
111 scale, we leveraged techniques from modern machine learning and computer vision, in  
112 particular the latent feature representations learned by large neural networks trained on

113 visual discrimination tasks (Radford et al., 2021; Simonyan & Zisserman, 2014), which  
114 have been shown in prior work to capture meaningful variation in human perceptual  
115 judgments about both natural images and drawings (Battleday, Peterson, & Griffiths,  
116 2020; Fan et al., 2018). We use these latent feature representations both to quantify the  
117 category-diagnostic variation in each drawing and to analyze the similarity structure in  
118 children’s unrecognizable drawings. In addition, we crowd-sourced part labels for each  
119 stroke in a subset of these drawings to quantify how the parts children included in their  
120 drawings changed across development. Finally, we administered drawing recognition  
121 tasks to measure how well children of different ages could identify which visual concept  
122 a given drawing was intended to convey.

123 This study makes a number of contributions relative to the prior literature. First,  
124 we collect, annotate, and share a large sample of children’s drawings from scribbles  
125 through sophisticated sketches, creating valuable resources for future research. Second,  
126 we develop an analytic approach suitable for exploring these drawings, which yields a  
127 number of intriguing findings around drawing development – including the presence of  
128 semantic information even in children’s unrecognizable drawings. Finally, we find  
129 evidence for the relation between developmental changes in children’s drawing abilities  
130 and their growing understanding of the visual concepts they are drawing. Older  
131 children include more diagnostic visual information and relevant object parts when  
132 producing line drawings, and these gains are not easily explainable by category  
133 exposure frequency or visuomotor development. Further, children’s developing ability to  
134 recognize drawings is related to the presence of category-diagnostic information in these  
135 drawings. Together, we provide a new set of tools and insights into the development of  
136 drawings and visual representations in childhood, which we hope will spur future  
137 research on this topic.



*Figure 1.* Top row: Museum kiosk where children participated, and examples of the tracing, drawing, and guessing trials. Bottom row: Example drawings from several categories; more red drawings contain more diagnostic visual features (as assessed by classifier evidence using VGG-19 FC6 features, see *Methods*).

138

## Results

### 139 Development of drawing production

140 A free-standing, child-friendly kiosk was installed at a local children's science  
 141 museum (see Figure 1, top row). Children used a touchscreen tablet attached to the  
 142 kiosk to produce their drawings. To evaluate how children's visuomotor abilities may  
 143 limit their ability to include the relevant visual features in their drawings (Freeman,  
 144 1987; Rehrig & Stromswold, 2018), we also included a set of shape-tracing trials in the  
 145 drawing production task to measure children's tracing skills (see Figure 1, top row).  
 146 After completing these tracing trials, children were verbally prompted to draw different

147 visual concepts. These categories were selected to include both animals and inanimate  
148 objects, as well as categories that are either commonly drawn (e.g., *cup*, *face*, *cat*) or  
149 less commonly drawn (e.g., *octopus*, *piano*, *camel*) by children (see *Methods* for more  
150 details, see Appendix A6). The final, filtered dataset contained 37,770 drawings of 48  
151 categories from  $N=8084$  children (average age: 5.33 years old; range: 2–10 years old;  
152 see *Appendix A* for detailed age demographics).

153 Measuring category-diagnostic information in such a large dataset of children's  
154 drawings spanning a wide variety of concepts poses a major analytical challenge. Until  
155 recently, researchers attempting to analyze even small drawing datasets had to develop  
156 *ad hoc* criteria for scoring drawings based on their intuitions about what the distinctive  
157 visual features could be for each target concept (e.g., handles for *mugs*) (e.g. Barrett &  
158 Light, 1976; Goodenough, 1963). Fortunately, recent advances in computer vision have  
159 made it possible to measure category-diagnostic information in images at scale by  
160 leveraging latent feature representations learned by large neural networks (Radford et  
161 al., 2021; Simonyan & Zisserman, 2014), although at some cost to interpretability, as  
162 these learned features are not guaranteed to map onto nameable object parts (e.g.,  
163 "handles"). Informed by this context, we use two approaches with complementary  
164 strengths to analyze our drawing dataset: first, we use model classifications to estimate  
165 the amount of category-diagnostic information in each drawing; second, we use  
166 crowdsourcing to identify which parts children included in their drawings in a general  
167 and scalable way.

168 Our first approach leverages the latent feature representations learned by large  
169 neural-network models to derive measures of a drawing's *recognizability* — how much  
170 category-diagnostic information it contains. Specifically, we analyzed the degree to  
171 which high-level visual features of each drawing could be used to decode the category  
172 that children intended to draw (e.g., *dog*), using features from VGG-19, a deep  
173 convolutional neural network pre-trained on Imagenet classification (Simonyan &  
174 Zisserman, 2014). Activations for each sketch were taken from the second-to-last layer of  
175 this model as prior work has shown that activations from deeper layers of convolutional

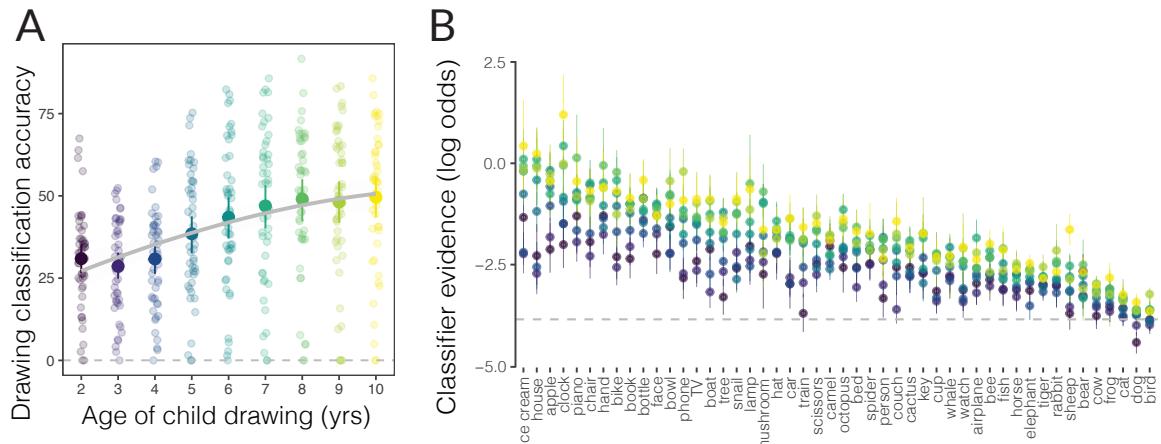
176 neural networks with a similar architecture correspond to the visual features that enable  
177 basic-level recognition (e.g., *cat* vs. *dog*) in both sketches and photographs (Fan et al.,  
178 2018; Yamins et al., 2014). These features were used to train logistic-regression  
179 classifiers to predict which of the 48 categories children were asked to draw (e.g., *couch*,  
180 *chair*) for sets of held-out drawings (see *Methods*), balanced across categories. For every  
181 drawing, this procedure thus yielded: (1) a binary classification score, indicating  
182 whether a given drawing contained the visual features that enabled basic-level  
183 recognition; and (2) a probability score for each of the 48 categories, capturing the  
184 degree to which a given drawing contained the visual features relevant to that category.

185 We then validated these VGG-19 model classifications by using embeddings from  
186 a modern contrastive language-image pre-training model (CLIP, Radford et al. 2021),  
187 which jointly trains an image and a text encoder to predict image and text pairings.

188 While less work has related the embeddings of this model class to either human  
189 behavioral or neural representations, CLIP outperforms other deep CNNs at recognizing  
190 visual concepts across different visual formats (Radford et al., 2021) and recent work  
191 suggests that CLIP embeddings show equal or better performance in predicting ventral  
192 stream responses (Conwell, Prince, Kay, Alvarez, & Konkle, 2022). CLIP classifications  
193 were obtained by assessing the similarity between model embeddings for each sketch to  
194 each category label, as in Radford et al., 2021. This method thus also yields both binary  
195 classification scores and probability scores for each of the 48 categories in the dataset.

196 As our second approach, we used a crowd-sourcing method to recruit human  
197 annotators to tag every stroke in a subset of  $N=2160$  drawings with a part label (see  
198 *Methods*). To ensure that these drawings were representative of the larger dataset, we  
199 chose 16 visual concepts (half animate, half inanimate, see *Methods*) and randomly  
200 sampled drawings from children 4-8 years of age. Using these annotations, we then  
201 analyzed changes in which parts children drew and how much they emphasized those  
202 parts in their drawings. The goal of these additional analyses was to provide insight  
203 into which specific elements within children's drawings change across development,  
204 giving rise to any changes in category-diagnostic information measured using model

205 classifications.



*Figure 2.* A. Proportion of drawings recognized as a function of children’s age; each dot represents the proportion of drawings that were correctly classified in a given category; the grey chance line represents 1/48 (number of categories in the dataset). B: The y-axis represents the log-odds probabilities (i.e. classifier evidence), binned by the age of the child who produced the drawing. Categories on the x-axis are ordered by average log odds probabilities for each category in descending order. Error bars represent bootstrapped 95% confidence intervals in both plots.

206           **Drawings of visual concepts become more recognizable across**

207       **childhood.** We found that children’s drawings increased in recognizability steadily  
 208 with age, as measured using model classification performance (VGG-19, Figure 2A,  
 209 Table 1; see validation using CLIP in Appendix, Figure A2, Table A3).<sup>1</sup> In an  
 210 additional study, we replicated this finding in a separate controlled experiment in which  
 211 a researcher was present while children produced their drawings ( $N=121$  children, ages  
 212 4-9 years), suggesting that the developmental changes we measured were not an artifact  
 213 of data collection at the museum kiosk (this is a subset of the data published in  
 214 BLINDED (in press), see Appendix, Figure A4, Table A4).

<sup>1</sup> We found that using features from deeper layers of VGG-19, rather than earlier layers, was critical to recovering these age-related changes, suggesting that drawings produced by older children primarily differed from those by younger children with respect to mid- and high-level visual features (see Appendix, Figure A5)

215       **More frequently drawn categories are not more recognizable.** What  
216   explains these gradual increases in recognizability? One way of accounting for these  
217   age-related differences is to suppose that younger children have had less practice  
218   drawing and are thus less well equipped to express what they know using this medium,  
219   despite having achieved a mature understanding of these visual categories. This account  
220   predicts that changes in recognizability are primarily driven by children's drawing  
221   experience with specific categories, either on their own or with caregivers and peers in  
222   educational contexts. If so, then frequently drawn categories (e.g., *trees*, *people*, *dogs*)  
223   should show the strongest developmental trends. To test this possibility, we asked  
224   parents to report how often their child produces drawings of each category ( $N=50$   
225   parents of children aged 3-10 years, *Methods*), revealing substantial variation in the  
226   frequency with which children tend to draw each of the categories in our stimulus set  
227   (see Appendix, Figure A6). We found converging evidence that drawings of more  
228   frequently practiced categories were no more recognizable and were not associated with  
229   stronger developmental trends; there was neither a main effect of drawing frequency on  
230   classification accuracy nor an interaction with age in a generalized linear mixed-effects  
231   model (see Table 1). This result was robust to the choice of model (VGG-19 and CLIP,  
232   see Appendix A3) and held when using human recognition scores in a separate  
233   controlled experiment (see Appendix, Figure A4, Appendix Table A4). Instead, we saw  
234   that many infrequently drawn categories (e.g. *ice cream*) had relatively high  
235   classification accuracy, while some frequently drawn categories (e.g. *dog*) had relatively  
236   low classification accuracy and were more likely to be confused with other similar  
237   categories (e.g., other animals) (see Figure 2B).

238       Figure 2B shows these developmental trends broken down by the category that  
239   children were intending to draw, highlighting the large amounts of variability across  
240   categories (see Appendix, A3 for validation using CLIP embeddings). We additionally  
241   examined whether other measures of frequency of experience in children's daily life  
242   might predict this item variation—for example, frequency in child-directed speech or  
243   all English-language books (see Appendix A5). However, we again found no discernable

244 relationship between these measures of frequency and the recognizability of children's  
245 drawings.

246 **Visuomotor control explains some but not all of changes in drawing**

247 **recognizability.** We anticipated that the recognizability of children's drawings would  
248 vary with their ability to control and plan their motor movements. Children spend  
249 countless hours across childhood both learning to write and practicing how to produce  
250 different shapes. Further, children's engagement with this drawing task could also  
251 reasonably vary as a function of age, with more skilled children spending more time,  
252 ink, or strokes on their drawings. We therefore measured the amount of time and effort  
253 children put into their drawings, and estimated children's visuomotor control via the  
254 simple shape tracing assessment task at the drawing kiosk. Children were instructed to  
255 trace both a relatively easy shape (a square) as well as a complex, novel shape that  
256 contained both curved and sharp segments (see Figure 1). For each participant, we used  
257 their performance on these two tracing trials to derive estimates of their tracing ability.  
258 Specifically, we obtained ratings of tracing accuracy from independent adult judges for a  
259 subset of tracings and then used these ratings to adapt an image registration algorithm  
260 (Sandkühler, Jud, Andermatt, & Cattin, 2018) to predict tracing scores for held-out  
261 tracings produced by children (see *Methods*). We found that tracing scores produced by  
262 the same participant were moderately correlated ( $r = .60$ ,  $t = 61.93$ ,  $df = 6754$ ,  $p <$   
263  $.001$ ,  $N = 6,756$ ), despite the irregular shape being harder to trace than the square.  
264 Thus, despite the brevity of this tracing assessment, the resulting measure had  
265 moderate reliability.

266 If age-related changes in drawing recognizability primarily reflect changes in  
267 visuomotor control (Freeman, 1987), then accounting for these more direct measures of  
268 visuomotor control ought to explain away the age-related variance we have observed so  
269 far. However, we still observed a robust main effect of age even after accounting for  
270 tracing abilities (Table 1) and other effort covariates (see Table 1), including the amount  
271 of time children spent drawing, the number of strokes in their drawings, and the amount  
272 of "ink" that they used (see *Methods*); this effect was robust to model choice (see

273 Appendix A3). These findings suggest that even though children's ability to control and  
 274 plan their motor movements does predict their ability to produce recognizable drawings,  
 275 this factor alone does not fully account for the observed developmental changes.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.690	0.173	-3.979	<0.001
Age (in years)	0.251	0.020	12.805	<0.001
Est. drawing frequency	-0.062	0.173	-0.356	0.721
Average tracing score	0.267	0.020	13.529	<0.001
Time spent drawing	0.039	0.021	1.868	0.062
'Ink' used	-0.031	0.020	-1.546	0.122
Number of strokes	0.008	0.018	0.477	0.634
Age*drawing frequency	0.017	0.017	1.042	0.298

Table 1

Model coefficients of a GLMM predicting the recognizability of each drawing (i.e. binary classification scores), including random intercepts for each category and participant. All predictors were z-scored so that coefficients are comparable.

### 276 Recognizable drawings contain more category-diagnostic information

277 **across development.** The above results suggest that children gradually improve  
 278 their ability to include diagnostic visual information in their drawings. However, they  
 279 are also consistent with an alternative account where younger children are just as able  
 280 to produce recognizable drawings when they are engaged with the task, but are less  
 281 likely to stay on task than older children and thus produce unrecognizable drawings  
 282 more often. To tease these two possibilities apart, we compared how much diagnostic  
 283 visual information was contained in drawings that were correctly classified.

284 For example, among drawings that were correctly recognized as *clocks*, did older  
 285 children also include visual information that more clearly set them apart from other  
 286 similar categories – for example, *watches*? Insofar as age-related improvements in  
 287 classification accuracy primarily reflect a decrease in the proportion of unrecognizable

288 drawings – rather than an increase in the quality of their recognizable drawings – we  
289 should expect drawings that were correctly classified to contain similar amounts of  
290 diagnostic visual information, regardless of whether they were produced by younger or  
291 older children.

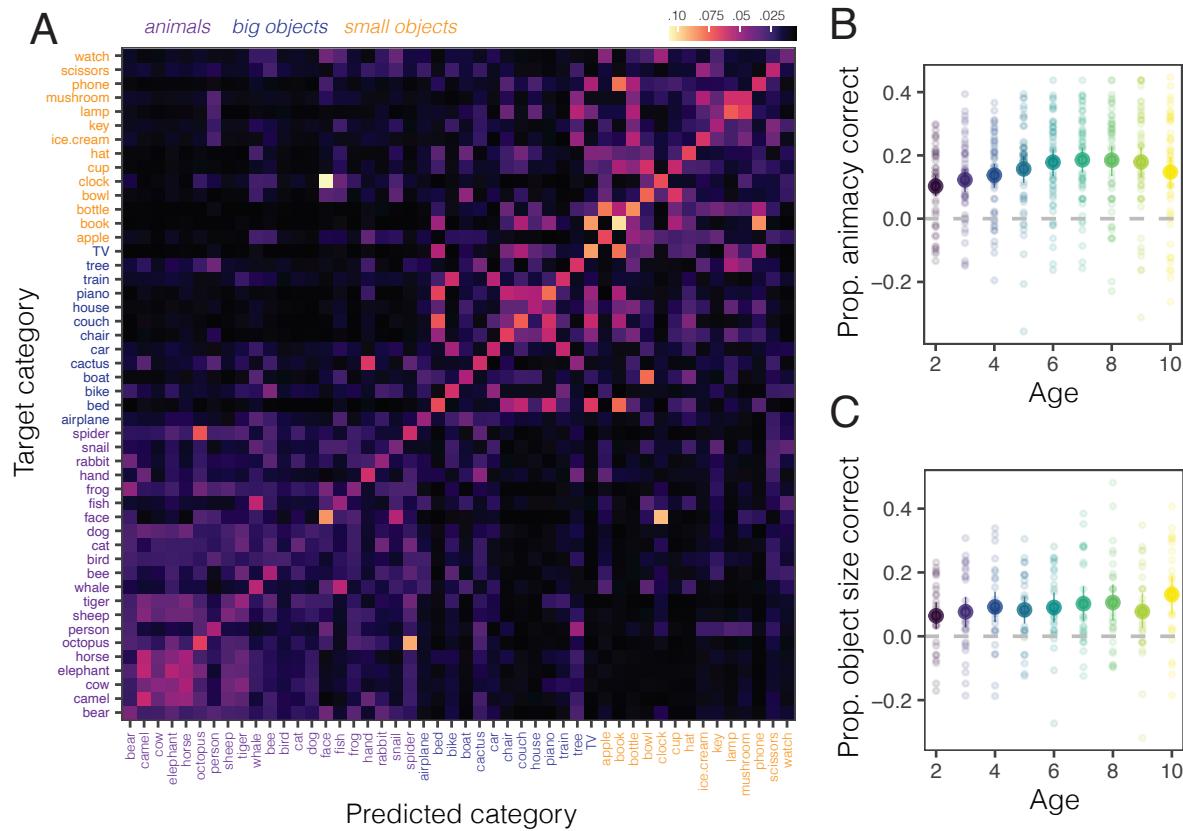
292 We found that even for drawings that were correctly classified (38.6% of the  
293 balanced subset of drawings,  $N=8,590$ ), the amount of diagnostic information they  
294 contained increased as a function of age, as measured by the log-odds probability  
295 assigned by the logistic-regression classifier to the target category (see *Methods*,  $B =$   
296  $0.111$ ,  $SE = 0.015$ ,  $df = 3544.18$ ,  $t = 7.354$ ,  $P < 0.001$ ) (see Appendix, Table A6. This  
297 analysis provides converging evidence that age-related improvements in children's  
298 abilities to produce recognizable drawings reflect a gradual increase in the amount of  
299 category-diagnostic information in their drawings.

300 **Unrecognizable drawings still contain semantically relevant**

301 **information.** Even if a child does not know the diagnostic features of *giraffes* or  
302 *rabbits*, they likely know that both are animals with four legs. Thus, this kind of coarse  
303 semantic information may still be contained in children's unrecognizable drawings.

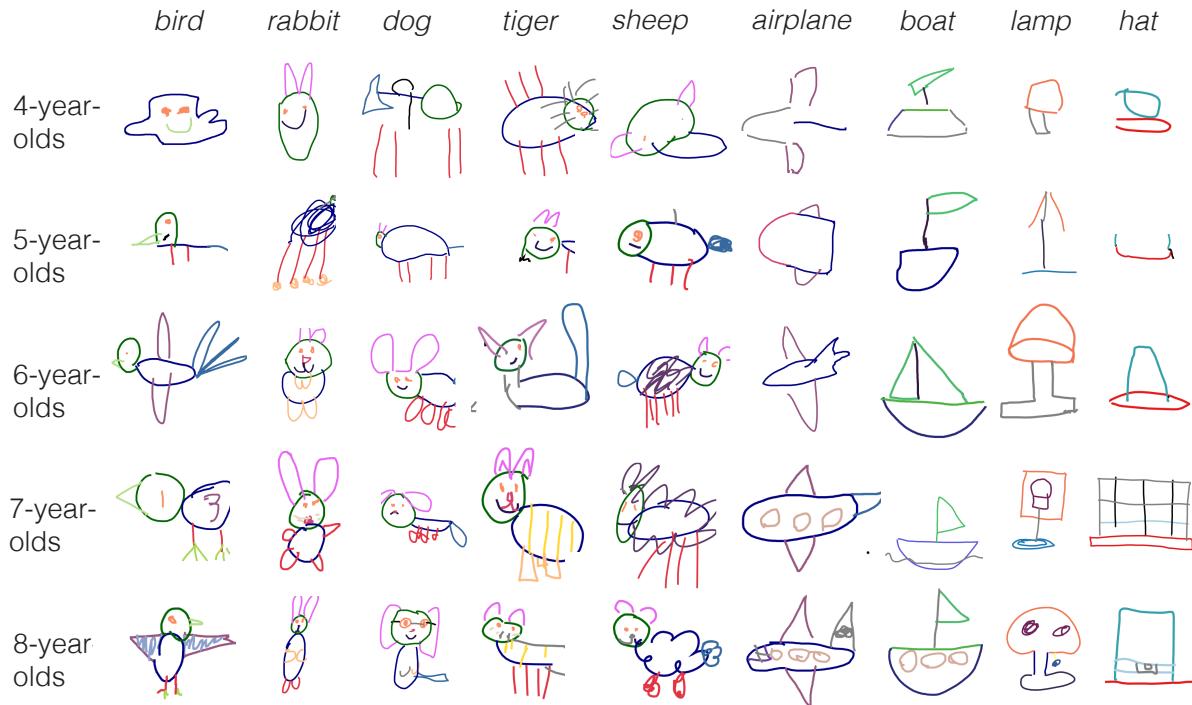
304 Indeed, prior work suggests that basic-level recognition – recognizing something as a  
305 *rabbit* – is not a pre-requisite for inferring semantic information. For example, adults  
306 can reliably judge the animacy (animate vs. inanimate) and real-world size of  
307 unrecognizable, textured images by inferring that animals tend to have high curvature  
308 and that larger, inanimate objects (e.g., couches) tend to have boxier shape structures  
309 (Long, Konkle, Cohen, & Alvarez, 2016; Long, Störmer, & Alvarez, 2017; Long, Yu, &  
310 Konkle, 2018) and children appear to be sensitive to these cues by the preschool years  
311 (Long, Moher, Carey, & Konkle, 2019).

312 Following this idea, we reasoned that even young children's children's misclassified  
313 drawings might contain information about the animacy and real-world size of the  
314 category they were intending to draw (see Figure 3A). To examine this possibility, we  
315 analyzed the patterns of misclassifications made by the logistic regressions and found  
316 that misclassified drawings reliably carried information about both real-world size (see



*Figure 3.* A: Classifier probabilities for the subset of drawings that were misclassified on the basis of VGG-19 embeddings (FC6). The y-axis shows the category children were intending to draw; the x-axis shows all of the categories in the dataset. Lighter values represent greater classifier probabilities assigned to a given category (see colorbar). (B,C). Proportion of misclassified drawings that contained the correct animacy/object size information of the target category (relative to baseline in the dataset). Each dot represents the proportion of drawings in a given category that had correct animacy/real-world size information relative to baseline at each age, respectively. Error bars represent bootstrapped 95% confidence intervals.

317 Figure 3C) and animacy (see Figure 3B) across all ages. We also found substantial  
 318 structure in the pattern of probabilities assigned by the classifier to the other categories  
 319 (see Figure 3A): for example, unrecognizable drawings of *an octopus* were often assigned  
 320 a high classifier probability for a *a spider*. These results suggest that children's  
 321 unrecognizable drawings are far from meaningless scribbles; instead, they contain  
 322 relevant semantic information about the category children were intending to draw.



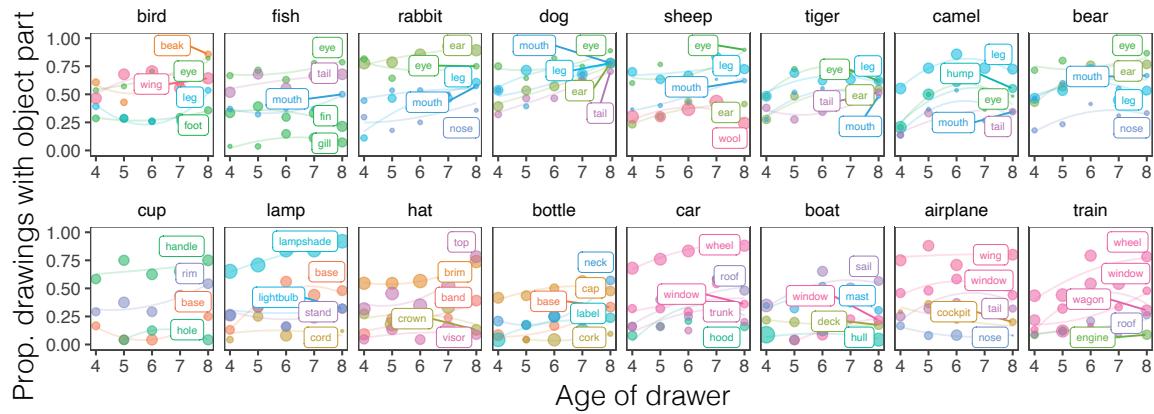
*Figure 4.* Example drawings from 4–8 year old’s, with part annotations. Each color represents object parts labels agreed upon by human annotators; grey lines represent strokes with multiple parts, and black lines represent unintelligible strokes.

**Drawings contain different semantic parts across development.** While

these findings so far provide strong evidence for age-related gains in the ability to produce recognizable drawings, it is not clear exactly what aspects of these drawings account for this improvement. A natural possibility is that children gradually learn which object parts to include and how much to emphasize those parts (e.g., long *ears* for rabbits) in their drawings (Tversky, 1989; Tversky & Hemenway, 1984) as they learn about the semantic properties of those categories. To explore this idea, we collected semantic part annotations of the visible object parts in a subset of  $N=2021$  drawings, and examined developmental changes in which parts children prioritized in their drawings throughout development.

Consistent with this idea, we found that drawings produced by older children

generally contained more unique semantic parts than drawings by younger children ( $B = 0.395$ ,  $SE = 0.041$ ,  $df = 2071$ ,  $P < .001$ ; Figure 5), lending support to the notion



*Figure 5.* Proportion of drawings that included a given object part for each object category as a function of children's age (see *Methods*). The size of each dot reflects the average emphasis (proportion of stroke length relative to the entire drawing) for each object part within each bin (max plotted part emphasis = .5); the top five most frequent object parts are included for each category excluding generic "body/head" parts.

336 that, across development, children learn to more effectively express what they know  
 337 about various visual concepts by enriching their drawings with additional semantic part  
 338 information (see Appendix, Figure A8). For some concepts, these gains appeared to be  
 339 specific to single part: for example, older children were more likely to produce *cups* with  
 340 *handles* and *cars* with recognizable *wheels*. For other concepts, however, age-related  
 341 changes were more complex: for example, in drawings of *bears*, both *ears* and *eyes*  
 342 appeared to change in prevalence and emphasis. Further, while most younger children's  
 343 drawings of *rabbits* included recognizable *ears*, which are in principle informative about  
 344 the category, many of them were still not recognizable as *rabbits*. Taken together, these  
 345 exploratory findings suggest that while there are clear age-related changes in the part  
 346 complexity of children's drawings, the mere presence of – or amount of emphasis on –  
 347 any particular part may not be sufficient to account for developmental variation in its  
 348 recognizability (see Appendix Figure A9). Instead, they suggest that visuospatial  
 349 information about what these parts look like and how they are arranged may be needed  
 350 to explain why drawings by older children are more recognizable than those by younger  
 351 ones. For example, *ears* on *rabbits* may need to be more elongated relative to the *head*

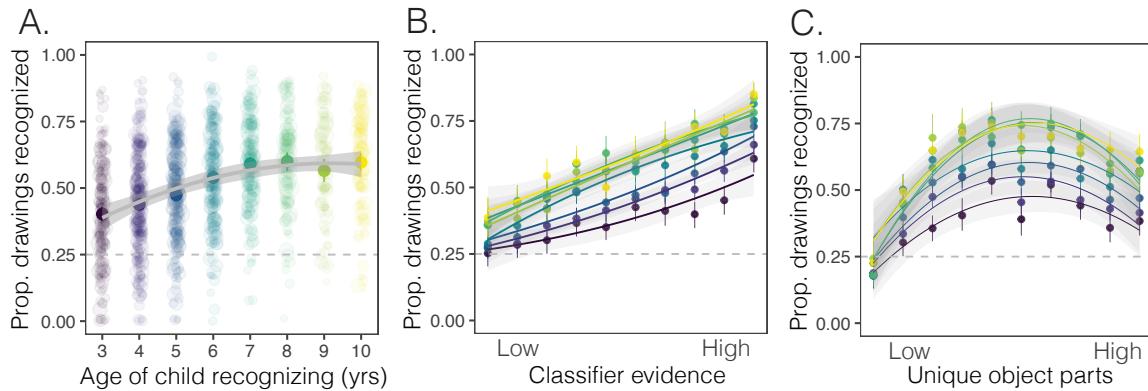
352 to provide a strong enough cue to category membership.

353 **Development of drawing recognition**

354 Why do children include more diagnostic visual information in their drawings as  
355 they grow older? One source of these developmental changes may be refinements in  
356 children's internal visual concepts. As children acquire more knowledge about the visual  
357 distinctions between visual concepts, children might more clearly represent the visual  
358 information that best distinguishes depictions of *rabbits* from *dogs*, for example, and  
359 may be able to use this information when recognizing drawings. This account thus  
360 predicts that children should improve over development in their ability to exploit visual  
361 information in drawings to recognize their intended meaning.

362 To explore this idea, we installed a “guessing game” in the same kiosk at the local  
363 science museum (see Figure 1, top right) where children guessed the category that an  
364 earlier child’s drawing referred to. These drawings were randomly sampled from the  
365 larger drawing dataset and thus varied in the degree to which they were recognizable  
366 and hence amount of diagnostic visual information they contained. This design choice  
367 allows us to examine how children’s visual recognition abilities vary when drawings  
368 contain differing amounts of diagnostic visual information.

369 Our goal in designing the visual recognition task was for it to be challenging yet  
370 not demand that children track a large number of comparisons. At the beginning of  
371 each session, children completed four practice trials in which they were cued with a  
372 photograph and asked to “tap the [vehicle/animal/object] that goes with the picture,”  
373 choosing from an array of four photographs of different visual concepts (see Figure 1).  
374 Children were then cued with drawings of these categories and responded using the  
375 same photograph buttons; photograph matching trials were also interspersed  
376 throughout the session as attention checks. We sequentially deployed four different  
377 versions of this task, including a different set of four perceptually similar categories in  
378 each (e.g., *hat*, *bottle*, *cup*, *lamp*). After exclusions, our dataset from this task included  
379 1,789 children ages 3–10 years (see *Methods*).



*Figure 6.* A. Drawing recognition as a function of the age of the child who participated in the guessing game; each dot represents data from one child who participated and is scaled by the number of trials they completed. (B,C). Drawing recognition data plotted separately by the age of the child participating as a function of the (B) amount of diagnostic visual information in each drawing, operationalized as the the *classifier evidence* assigned to each sketch relative to the distractor categories and (C) the number of unique object parts in each drawing. Both variables are binned into deciles for visualization purposes. Error bars represent bootstrapped 95% confidence intervals.

380 Overall, we found that children became steadily better at identifying the category  
 381 that a drawing referred to (see Figure 6A). In contrast, performance on photograph  
 382 matching trials was relatively similar across ages. All children whose data were included  
 383 in our analyses scored >75% correct on photograph trials and average accuracy in each  
 384 group ranged from  $M=90\text{-}93\%$  correct. Thus, variation in drawing recognition accuracy  
 385 is unlikely to be explained by generic differences in motivation or task engagement.

386 **Children's drawing recognition improves over development.** We next  
 387 evaluated the idea that children's ability to exploit category-diagnostic visual  
 388 information during recognition improves over childhood. We first tested how children's  
 389 drawing recognition ability varied with respect to the amount of diagnostic visual  
 390 information in a given drawing. To do so, for each drawing that appeared in the  
 391 guessing games, we measured diagnostic visual information. We fit a 4-way logistic  
 392 regression classifier trained on the VGG-19 features extracted from the drawings  
 393 presented in each guessing game (see *Methods*) and measured diagnostic information as

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	0.050	0.121	0.412	0.680
Classifier evidence	0.477	0.046	10.405	<0.001
Recognizer Age	0.317	0.019	16.777	<0.001
Classifier evidence*Recognizer Age	0.062	0.014	4.246	<0.001

Table 2

*Model coefficients of a GLMM predicting visual recognition performance as a function of recognizer age and classifier evidence. All predictors were z-scored such that coefficients are comparable.*

394 the log-odds ratio between the intended category and the foil categories. That is, the  
 395 diagnostic information for a *dog* drawing was defined relative to its perceptual similarity  
 396 to the other choices in the recognition task (i.e., *bird*, *fish*, *rabbit*). We then fit a  
 397 generalized linear mixed effects models predicting children’s recognition performance  
 398 with child’s age, this metric of diagnostic visual features, and their interaction as fixed  
 399 effects (see *Methods* for further details and robustness checks using CLIP).

400 Drawings with more diagnostic visual information were better recognized across  
 401 all ages (see Table 2, Figure 6B, CLIP robustness check in Appendix, Table B2). Yet  
 402 older children were also better able to capitalize on graded differences in the diagnostic  
 403 visual information in drawings when recognizing them (see Figure 6B), evidenced by an  
 404 interaction between classifier evidence and recognizer age in both cases. This result held  
 405 when we restricted our analyses to a subset of high-performing children who performed  
 406 at ceiling on photograph matching trials (see *Appendix*, B1) suggesting that these effects  
 407 are unlikely to be driven by a differences across development in either engagement or in  
 408 the ability to match drawings with the picture-cue buttons (see individual category  
 409 effects in *Appendix*, B3). Children became steadily better over development at using  
 410 diagnostic visual information to recognize the intended meaning of line drawings.

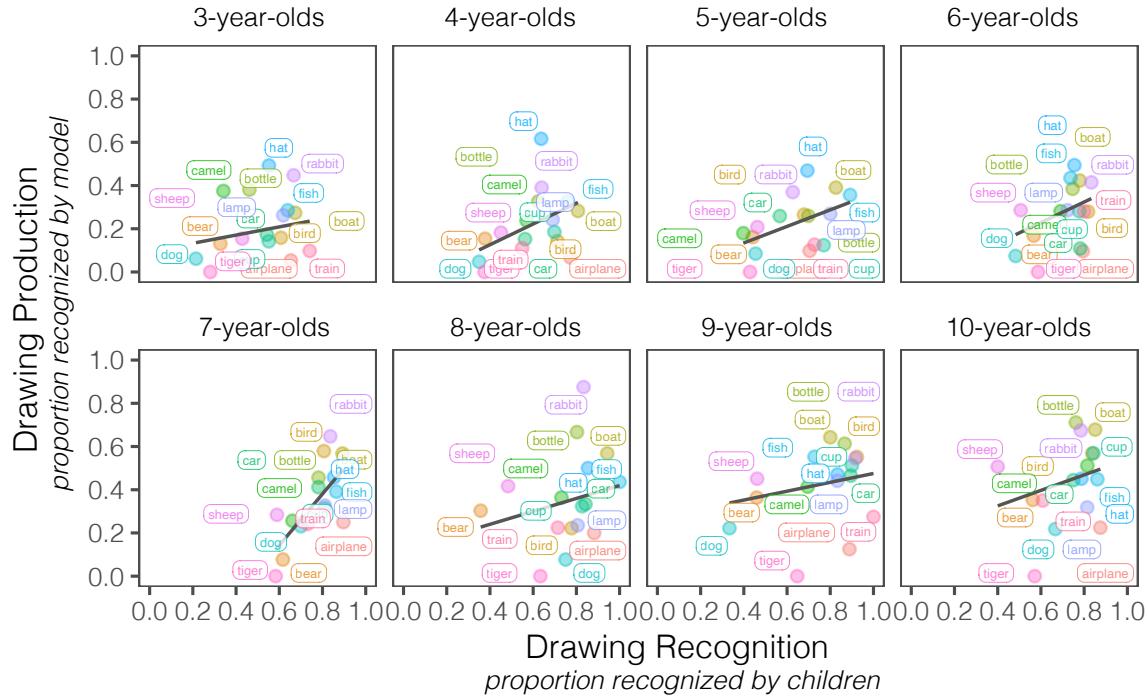
411 **Children’s drawing recognition varies with part information.** Finally, we  
 412 evaluated how children’s abilities to use object part information during visual

413 recognition changes across development (Juttner et al., 2016; Mash, 2006). Specifically,  
414 we examined how children’s recognition accuracy varied with the number of unique  
415 recognizable object parts included in the drawings they tried to recognize. We again fit  
416 a generalized linear mixed effect model to the recognition data, modeling the interaction  
417 between the number of unique parts in each drawing and the age of the child recognizing  
418 the drawing, with the maximal random effects structure supported by our data.

419 We found that drawings with more unique object parts tended to be better  
420 recognized – but, unexpectedly, that drawings with many object parts were less well  
421 recognized than drawings with an intermediate number of parts (see Figure 6C) though  
422 there was substantial variation across categories (see Appendix, Figure A1, see  
423 Table B3). Critically, we again found that older children were better able to capitalize  
424 on increasing object part information during recognition, as evidenced by a significant  
425 interaction between the number of unique parts in the drawings and the age of  
426 child recognizing the drawings. In other words, children’s ability to integrate additional  
427 object part information during recognition changed across development (see all  
428 coefficients in Appendix, Table B3). These additional analyses using object part  
429 annotations suggest that children’s ability to use diagnostic visual information during  
430 recognition matures gradually throughout childhood.

### 431 Relationship between visual production and recognition

432 So far, our descriptive results suggest that both visual production and recognition  
433 of drawings develop gradually and in parallel throughout childhood. To what degree do  
434 these developmental trajectories reflect changes in the same mental representations of  
435 visual concepts? Insofar as children’s abilities to recognize what drawings mean and to  
436 produce meaningful drawings both rely on a shared mental representation, then their  
437 ability to produce a drawing of a *dog* may be correlated with their ability to recognize a  
438 drawing of a *dog*, as in adults (Perdreau & Cavanagh, 2014). To explore this possibility,  
439 we borrow an analysis technique used in language acquisition (Braginsky, Yurovsky,  
440 Marchman, & Frank, 2019), where variation in word production is well-predicted by



*Figure 7.* Each dot represents a category (e.g., hat) at a given age (in years), where the y-axis value represents how well children of that age produced recognizable drawings of that category (as assessed by CLIP model classifications) and the x-axis value represents how well children of that age were able to *recognize* the top 30% most recognizable drawings of that category (as assessed by accuracy in the 4AFC recognition games). Independent sets of drawings are analyzed in each case.

<sup>441</sup> independent data about that word, e.g., the frequency of a word in English books

<sup>442</sup> (Goodman, Dale, & Li, 2008).

<sup>443</sup> We thus explored how well variation in visual production is related to visual  
<sup>444</sup> recognition at the category level, acknowledging the exploratory nature of these  
<sup>445</sup> analyses.<sup>2</sup> To do so, we examined children's visual production and recognition abilities  
<sup>446</sup> using independent sets of drawings of the same categories. To estimate drawing  
<sup>447</sup> recognition ability, we used children's performance on the guessing games to calculate

<sup>2</sup> While some children may have both contributed drawings and participated in the recognition game with different categories, these sessions were anonymous, and thus we do not have access to within-child data.

448 how often children of a given age, on average, were able to recognize drawings of a given  
449 category. To ensure that we were examining children's visual recognition abilities for  
450 relatively recognizable drawings, we analyzed how well children could recognize the top  
451 30 percent most recognizable drawings of each category (using CLIP classification  
452 probabilities, see *Methods*). To estimate drawing production abilities, we calculated how  
453 often children at each age could produce recognizable drawings of a given category (e.g.,  
454 *a dog*), as assessed by CLIP model recognition scores. For these exploratory analyses,  
455 we used CLIP model classifications, as CLIP showed less dramatic category variation  
456 relative to VGG-19 classifications (see Appendix, Figure A2). In addition, within the  
457 independent set of sketches used to assess children's recognition, CLIP showed a higher  
458 correlation with children's recognition behaviors (aggregating across individual sketches;  
459 VGG-19,  $r=.28$ ; CLIP,  $r=.43$ ).

460 Overall, we found that children's visual production and visual recognition abilities  
461 were positively related at the category-level at all ages (see Figure 7). For example,  
462 while *dogs* and *sheep* were consistently both harder to produce and to recognize, *rabbits*  
463 and *hats* were easier to produce and recognize. Thus, these exploratory results suggest  
464 relative consistency across categories in these two tasks, suggesting that children's  
465 ability to perform well in both tasks may rely on a shared visual representation, and  
466 paving the way for future work that seeks to understand the sources of this category  
467 variation using within-child, controlled experiments.

## 468 General Discussion

469 We conducted a systematic investigation of how children produce and recognize  
470 line drawings of a wide range of everyday visual concepts across development (2-10  
471 years of age). We developed a large dataset of children's digital drawings (>37K) and  
472 capitalized on innovations in machine learning to quantify changes in children's  
473 drawings across development. We found robust improvements in children's ability to  
474 include diagnostic visual information via recognizable object parts in their drawings,  
475 and these developmental changes were not reducible to either increased effort or better

476 visuomotor abilities. Further, we found that children's unrecognizable drawings still  
477 contained information about the animacy and real-world size of the visual concepts they  
478 were trying to depict; highlighting an intermediate stage between scribbles and fully  
479 recognizable drawings. We also found improvements throughout childhood in children's  
480 ability to recognize each other's drawings, particularly in their ability to capitalize on  
481 diagnostic visual information during drawing recognition tasks. Together, these results  
482 document parallel developmental changes in how children use diagnostic visual  
483 information when producing and recognizing freehand line drawings, suggesting that  
484 refinements in children's visual concepts may underlie improvements across both tasks.  
485 As children's perceptual abilities (Bova et al., 2007; Natu et al., 2016), semantic  
486 knowledge (Tversky, 1985; Vales, Stevens, & Fisher, 2020), and visuomotor skills (Li &  
487 James, 2016) evolve across childhood, children's changing visual concepts influence both  
488 how children produce and recognize freehand line drawings.

489 More broadly, the present work highlights how the combination of modern  
490 machine-learning methods and larger-scale datasets of naturalistic behaviors can  
491 contribute to theoretical progress in developmental science. By collecting rich data from  
492 many participants over a large developmental age range, we can more precisely estimate  
493 graded changes in children's abilities and the degree to which these trajectories vary  
494 across categories. In turn, our use of innovations in computer vision and computational  
495 modeling allow the analysis of the entirety of this large dataset, capturing variation  
496 across both unrecognizable and recognizable drawings in a single analytic approach  
497 (which would have been intractable with human ratings, e.g., how similar each *dog*  
498 drawing was to every other category in the dataset). Using this approach, we were able  
499 to distinguish variability in children's drawings due to a range of different  
500 developmental processes – including motor skill and task effort – from variability related  
501 to visual concept knowledge. We believe that this work paints a more accurate picture  
502 of developmental change and opens up new avenues for investigating the various factors  
503 that shape visual concepts throughout development (Karmiloff-Smith, 1990; Minsky &  
504 Papert, 1972) both using large-scale datasets and controlled, within-child experiments

505 that directly relate visual production and recognition and examine item variation.

506 **Variation across visual concepts**

507 Our exploratory analyses suggested that children's abilities to produce and  
508 recognize drawings were correlated at the category-level at each age, e.g., drawings of  
509 *dogs* were both harder to produce and to recognize. Further, estimates of how often  
510 these items tend to be drawn or experienced did not explain this variation. Why might  
511 some categories be easier to draw and to recognize?

512 One possibility is that these item effects are related to other metrics of visual  
513 experience with different categories beyond what we have measured: Perhaps exposure  
514 frequency in educational materials or children's media were not adequately captured by  
515 our surveys. Or perhaps children may experience more invariant exemplars of some  
516 categories, making it easier to identify and draw those categories. For example, while  
517 *dogs* may vary substantially, *ice cream cones* have a relatively more invariant form.  
518 Children may in turn develop more refined visual representations for more frequently  
519 experienced and more invariant items, leading to more recognizable drawings.

520 Any effects of exposure frequency or form variability might interact with the  
521 degree to which a category has a 3D shape structure that can be easily depicted using a  
522 two-dimensional line drawing. For example, canonical *mushrooms* have a relatively  
523 simple shape, whereas *rabbits* have many sub-parts that need to be depicted (i.e., legs,  
524 ears, nose, mouth, tail), and children may also struggle to arrange these sub-parts in a  
525 way that conveys the meaning of the visual concept (i.e., the correct ratio between the  
526 size of the rabbit's ears and head). In turn, these shape structures may lead different  
527 categories to have more or less typical iconic representations that children are accessing  
528 when producing and recognizing line drawings. For example, *trains* are often depicted  
529 as steam trains, as modern trains can be hard to distinguish from other vehicles as line  
530 drawings. Future work may be able to use our dataset to examine the relationship of  
531 these factors to visual production and recognition, as we purposefully included items  
532 that vary along these dimensions.

533

**534 From scribbles to categories**

535 The present work also highlights the gradual progression in children's drawings  
536 from exploratory scribbles through an intermediate stage (Morra & Panesi, 2017) where  
537 their drawings may not unequivocally convey a specific visual concept (e.g., a *giraffe*),  
538 while still containing enough visual information to be recognizable as an "animal." An  
539 intriguing implication is that the mechanisms by which children reliably produce such  
540 semantically ambiguous drawings might be related to mid-level visual representations in  
541 the brain that are sensitive to coarse distinctions between broad classes of visual objects  
542 (i.e., large inanimate vs. small inanimate vs. animals) without relying on category-level  
543 distinctions (Long et al., 2016, 2019, 2017). Moreover, drawing tasks may allow children  
544 to convey this kind of partial knowledge about a visual concept that may otherwise be  
545 difficult for them to express verbally. And as children learn more about a specific visual  
546 concept – for example, that *giraffes* have longer legs than *antelopes* or *elephants* – these  
547 incremental gains in conceptual knowledge may manifest in their drawings, even if they  
548 are not yet clearly recognizable as a *giraffe*. This work thus showcases the potential of  
549 drawing production tasks for examining graded changes in how children's knowledge  
550 about visual concepts grows and changes across development.

**551 Possible learning mechanisms**

552 Several learning mechanisms are consistent with the developmental changes we  
553 observed. One possibility is that children become better visual communicators as they  
554 learn which visual features are most effective at conveying category membership  
555 through the process of producing drawings. In turn, this process of using drawings and  
556 other visual modalities to communicate various visual concepts may have downstream  
557 effects on children's ability to recognize drawings of them. Indeed, both drawing experts  
558 (Perdreau & Cavanagh, 2013b) and naïve adults who practice drawing similar categories  
559 (Fan et al., 2020) show better enhanced visual recognition abilities of these categories.  
560 Such a mechanism would be consistent with prior work suggesting that learning to

561 produce letters by hand can support subsequent letter recognition (James, 2017;  
562 Longcamp, Zerbato-Poudou, & Velay, 2005), with recent findings pointing towards the  
563 variability of visual forms seen while learning to write as a key factor (Li & James,  
564 2016). Thus, the process of iteratively producing and recognizing drawings of visual  
565 concepts could cause these parallel developmental changes in both domains. Contra a  
566 strong version of this account, however, we did not find strong effects of drawing  
567 practice at the category level in the present data: for example, *ice cream cones* were  
568 among the best recognized categories and estimated (by parents) to be among the least  
569 practiced by children. In addition, we did not see any obvious changes in the  
570 developmental trajectory around 6 years of age when most children start to write  
571 (though this achievement might cause smaller changes that we could not detect in our  
572 data); rather, we observed evidence for gradual growth throughout the entire age range.

573 A second, non-exclusive possibility is that children are explicitly learning the  
574 diagnostic features of categories as they enrich their semantic knowledge. For example,  
575 children may learn about the functional properties of different attributes: *camels* have  
576 humps to store water, *clocks* have numbers to tell time, and *whales* spout water because  
577 they need to breathe. As a result, children may come to more clearly represent which  
578 visual features are diagnostic of different categories and why. In turn, this semantic  
579 knowledge could percolate into children's visual concepts and thus be accessed both  
580 when children draw an object and when they recognize it. This possibility aligns with a  
581 wealth of evidence suggesting that continual learning about different categories  
582 throughout the early school years shapes children's categorization abilities. For  
583 example, children change in how they think about the diagnosticity of different  
584 semantic properties across development: in early childhood, the fastest cheetah – that  
585 is, the exemplar with the most extreme value on some property – tends to be seen as  
586 the best and the most representative cheetah (Foster-Hanson & Rhodes, 2019). At the  
587 same time, taxonomic groupings become increasingly important both in children's  
588 explicit conceptual judgements (Tversky, 1985) and when children spontaneously  
589 arrange different visual concepts (e.g., wild vs. farm animals, Vales et al., 2020). Thus,

590 children's evolving semantic knowledge could shape the visual features children use both  
591 when producing and recognizing different visual concepts.

592 A third possibility, again not mutually exclusive with the other two, is that  
593 children are implicitly learning category-diagnostic information through the process of  
594 visual categorization itself: through repetitively viewing and categorizing depictions,  
595 real-life examples, and photographs of these different categories. Indeed, the neural  
596 networks used here to categorize drawings did not have visuomotor experience drawing  
597 or training about the semantic properties of these categories. Thus in principle it is  
598 possible that children could be refining their visual concepts without substantial  
599 involvement from other cognitive or sensorimotor systems.

## 600 **Limitations and future directions**

601 There are various limitations to the generalizability of these findings that future  
602 work could address. First, while these datasets are large and sample heterogeneous  
603 populations, all drawings and recognition behaviors were collected at a single  
604 geographical location, limiting the generalizability of these results to children from  
605 other cultural or socioeconomic backgrounds (Henrich, Heine, & Norenzayan, 2010).  
606 Children in different contexts may spend considerably more or less time viewing and  
607 producing depictions of different categories, and thus could reasonably affect how  
608 how they represent them. Further, different cultural contexts have different conventions  
609 for depicting visual concepts, as evidenced by drawings from adults (Lewis,  
610 Balamurugan, Zheng, & Lupyán, 2021). However, there are likely to be some aspects of  
611 drawing production and interpretation that are broadly shared across cultural contexts  
612 (Cavanagh, 2005; Hertzmann, 2020), given prior work that has investigated picture  
613 comprehension in communities with modest exposure to Western visual media  
614 (Deregowski, 1989; Kennedy & Ross, 1975). Moreover, there is evidence from earlier  
615 work that some of this convergence may reflect evolutionarily conserved visual  
616 processing mechanisms, as non-human primates exhibit can recognize the  
617 correspondence between line drawings and their real-world referents (Itakura, 1994;

618 Tanaka, 2007). Future work that examines drawings across different cultural contexts in  
619 both adults (Lewis et al., 2021) and children will help quantify the consistency and  
620 variability in how we represent and depict visual concepts.

621 Second, while we imposed strong filtering requirements, we were not present while  
622 the children were drawing or guessing at the kiosk and thus cannot be sure that we  
623 eliminated all sources of noise or interference. Many sources of additional interference  
624 would only generate noise in our data, though, rather than creating specific age-related  
625 trends. Nonetheless, we replicated our main experimental results on drawing production  
626 in a controlled, experimental context with a smaller set of categories (see Appendix,  
627 Figure A4).

628 Third, since these datasets are cross-sectional, they do not directly relate visual  
629 production and recognition abilities at the individual level. Our exploratory,  
630 category-level analyses suggest variation in these two abilities are correlated across  
631 development; ultimately, however, within-child measurements will be necessary to  
632 confirm that changes in children's visual concepts underlie the observed changes in both  
633 tasks. In addition, these correlational analyses can only provide hints as to whether  
634 changes in visual production cause changes in visual recognition or vice versa.  
635 Finer-grain, within-child training studies (as in Bremner & Moore, 1984) could provide  
636 traction on the direction of causality between visual production and recognition.

## 637 Conclusion

638 Our results call for further systematic, experimental investigations into the kinds  
639 of experience – including visuomotor practice, semantic enrichment, and visual exposure  
640 – that may influence visual production and recognition in children, and we hope that  
641 the open datasets and tools we have created here will open up new avenues for such  
642 future work. We propose that a full understanding of how children produce and  
643 recognize drawings of visual concepts will allow a unique and novel perspective on the  
644 both the development and the nature of visual concepts: the representations that allow  
645 us to easily derive meaning from what we see.

646

## Methods & Materials

### 647 Drawing Station Details

648 While the interface was designed to be navigable by children, the first page of the  
649 drawing station showed a short consent form and asked parents to enter their child's age  
650 in years; no other demographic information was collected. Afterwards, video prompts of  
651 an experimenter guided the child through the rest of the experiment; an initial video  
652 stated that this game was "only for one person at a time" and asked children to "draw  
653 by themselves." Every session at the drawing station started with tracing trials before  
654 moving on to the category prompts ("What about a [*couch*]? Can you draw a [*couch*]?").  
655 Children could stop the experiment at any time by pressing a stop button; each trial  
656 ended after 30 seconds or after the child pressed the "next" button. Six different sets of  
657 eight category prompts rotated at the station, yielding drawings from a total of 48  
658 categories (see Appendix Figure B1, *airplane*, *apple*, *bear*, *bed*, *bee*, *bike*, *bird*, *boat*,  
659 *book*, *bottle*, *bowl*, *cactus*, *camel*, *car*, *cat*, *chair*, *clock*, *couch*, *cow*, *cup*, *dog*, *elephant*,  
660 *face*, *fish*, *frog*, *hand*, *hat*, *horse*, *house*, *ice cream*, *key*, *lamp*, *mushroom*, *octopus*,  
661 *person*, *phone*, *piano*, *rabbit*, *scissors*, *sheep*, *snail*, *spider*, *tiger*, *train*, *tree*, *TV*, *watch*,  
662 *whale*); these categories were also chosen to overlap with those in the QuickDraw  
663 database of adult drawings (<https://github.com/googlecreativelab/quickdraw-dataset>).  
664 Each set of category prompts that rotated at the station thus included both animate  
665 and inanimate categories as well as commonly and infrequently drawn categories;  
666 category prompts were presented in a random order.

### 667 Drawing Dataset Filtering & Descriptives

668 Given that we could not easily monitor all environmental variables at the drawing  
669 station that could impact task engagement (e.g., ambient noise, distraction from other  
670 museum visitors), we anticipated the need to develop robust and consistent procedures  
671 for data quality assurance. We thus adopted strict screening procedures to ensure that  
672 any age-related trends we observed were not due to differences in task compliance  
673 across age. Early on, we noticed an unusual degree of sophistication in 2-year-old

674 participants' drawings and suspected that adult caregivers accompanying these children  
675 may not have complied with task instructions to let children draw on their own. Thus,  
676 in subsequent versions of the drawing game, we surveyed participants to find out  
677 whether another child or an adult had also drawn during the session; all drawings where  
678 interference was reported were excluded from analyses. Out of 11797 subsequent  
679 sessions at the station, 3094 filled out the survey, and 719 reported interference, 6.09%  
680 of participants; these participants' drawings were not rendered or included in analysis.  
681 When observing participants interacting with the drawing station, we noted that most  
682 children's parents did not fill out the survey because they were either talking to other  
683 parents or taking care of a sibling. Further, while children could contribute drawing  
684 data more than once if they chose, this did not occur during our structured observation  
685 of the kiosk. This protocol was approved by both the Institutional Review Board at  
686 Stanford University (43992, Development of Children's Drawing Abilities).

687 Raw drawing data were then screened for task compliance using a combination of  
688 manual and automated procedures (i.e., excluding blank drawings, pure scribbles, and  
689 drawings containing words). A first subset of drawings ( $N = 15594$  drawings) was  
690 filtered manually by one of the authors, resulting in  $N = 13119$  drawings after  
691 exclusions (15.8% exclusion rate); subsequently, drawing filtering was crowd sourced via  
692 Prolific. 390 participants first completed a practice round demonstrating valid and  
693 invalid drawings and then viewed 24 drawings from a intended category at a time and  
694 selected the invalid drawings they judged to come from off-task participants.  
695 Participants were reminded that unrecognizable drawings were still "valid" drawings,  
696 and could proceed to the next category only after selecting a 'catch' invalid drawing.  
697 Each drawing in the dataset was viewed at least twice by two different participants. To  
698 be conservative, any drawing that was marked as 'invalid' by a participant was excluded  
699 from the dataset. These stringent filtering criteria resulted in the exclusion of an  
700 additional 9897 drawings, leading to an overall exclusion rate of 24.57% of the drawings  
701 and a final set of 37770 drawings from 8084 sessions. In the final dataset, there were  
702 more younger than older children, despite filtering; see *Appendix Table A1* for a

703 complete summary.

#### 704 Experimental Dataset Procedure

705 In a separate experiment, children were seated in front of a touchscreen tablet  
706 with a trained experimenter. As in the larger dataset, children completed two  
707 shape-tracing trials, and then children produced drawings of 12 familiar object  
708 categories (*airplane, bike, bird, car, cat, chair, cup, hat, house, rabbit, tree, watch*)  
709 which were randomly assigned to different cue-types (verbal vs. picture). In this paper,  
710 we analyze only verbal-cued drawings for sake of comparison to the drawing station  
711 dataset. 135 children participated in the experiment; 6 participants were excluded, (3)  
712 for skipping more than 6 drawing trials and (3) for scribbling three or more times in a  
713 row. Six additional participants were tested but their data was not recorded due to a  
714 technical error, and two participants never advanced past the practice trials, leading to  
715 a final sample of 121 children. No additional demographic data was recorded about the  
716 participants. This protocol was approved by the Institutional Review Board at Stanford  
717 University (43992, Development of Children's Drawing Abilities).

#### 718 Measuring Tracing Accuracy

719 We developed an automated procedure for evaluating how accurately participants  
720 performed the tracing task that was validated against empirical judgments of tracing  
721 quality. We decompose tracing accuracy into two components: a *shape error* component  
722 and a *spatial error* component. Shape error reflects how closely the participant's tracing  
723 matched the contours of the target shape; the spatial error reflects how closely the  
724 location, size, and orientation of the participant's tracing matched the target shape.

725 To compute these error components, we applied an image registration algorithm,  
726 AirLab (Sandkühler et al., 2018), to align each tracing to the target shape, yielding an  
727 affine transformation matrix that minimized the pixel-wise correlation distance between  
728 the aligned tracing,  $T$ , and the target shape,  $S$ :  $Loss_{NCC} = -\frac{\sum S \cdot T - \sum E(S)E(T)}{N \sum Var(S)Var(T)}$ , where  
729  $N$  is the number of pixels in both images.

730        The shape error was defined by the final correlation distance between the aligned  
731 tracing and the target shape. The spatial error was defined by the magnitude of three  
732 distinct error terms: location, orientation, and size error, derived by decomposing the  
733 affine transformation matrix above into translation, rotation, and scaling components,  
734 respectively. In sum, this procedure yielded four error values for each tracing: one value  
735 representing the shape error (i.e., the pixel-wise correlation distance) and three values  
736 representing the spatial error (i.e., magnitude of translation, rotation, scaling  
737 components).

738        Although we assumed that both shape and spatial error terms should contribute  
739 to our measure of tracing task performance, we did not know how much weight to  
740 assign to each component to best predict empirical judgments of tracing quality. In  
741 order to estimate these weights, we collected quality ratings from adult observers  
742 ( $N=70$ ) for 1325 tracings (i.e., 50-80 tracings per shape per age), each of which was  
743 rated 1-5 times. Raters were instructed to evaluate “how well the tracing matches the  
744 target shape and is aligned to the position of the target shape” on a 5-point scale.

745        We fit an ordinal regression mixed-effects model to predict these 5-point ratings,  
746 which contained correlation distance, translation, rotation, scaling, and shape identity  
747 (square vs. star) as predictors, with random intercepts for rater. This model yielded  
748 parameter estimates that could then be used to score each tracing in the dataset  
749 ( $N=14372$  tracings from 7612 children who completed at least one tracing trial). We  
750 averaged scores for both shapes within session to yield a single tracing score for each  
751 participant.

## 752 Measuring effort covariates

753        For each drawing trial, children had up to 30 seconds to complete their drawings  
754 with their fingers. We recorded both the final drawings and the parameters of each  
755 stroke produced by children at the drawing station, allowing us to estimate the amount  
756 of time children put into their drawings. As a second measure of effort, we also counted  
757 the number of strokes that children put into a given drawing. Finally, we estimated the

758 proportion of the drawing canvas that was filled (e.g., ‘ink used’) by computing the  
759 proportion of each final drawing that contained non-white pixels.

760 **Estimating drawing recognizability**

761 **VGG-19: Visual Encoder.** To encode the high-level visual features of each  
762 sketch, we used the VGG-19 architecture (Simonyan & Zisserman, 2014), a deep  
763 convolutional neural network pre-trained on Imagenet classification. For our main  
764 analysis, we used model activations in the second-to-last layer of this network, which is  
765 the first fully connected layer of the network (FC6), as prior work suggests that it  
766 contain more explicit representations of object identity than earlier layers (Fan et al.,  
767 2018; Long, Fan, & Frank, 2018; Yamins et al., 2014). Raw feature representations in  
768 this layer consist of flat 4096-dimensional vectors. to which we applied channel-wise  
769 normalization across all filtered drawings in the dataset. For additional analyses using  
770 the earlier convolutional layers, we first applied spatial averaging over the outputs of  
771 each layer to reduce their dimensionality, as in Fan et al., 2018, before also applying  
772 channel-wise normalization.

773 **VGG-19: Logistic regression classifiers.** Next, we used these features to  
774 train object category decoders. To avoid any bias due to imbalance in the distribution  
775 of drawings over categories (since groups of categories ran at the station for different  
776 times), we sampled such that there were an equal number of drawings of each of the 48  
777 categories ( $N=22,272$  drawings total). We then trained a 48-way logistic classifier with  
778 L2 regularization (tolerance = .1, regularization = .1), and used this classifier to  
779 estimate the category labels for a random held-out subset of 96 drawings (2 drawings  
780 from each category). No additional metadata about the age of the child who produced  
781 each sketch was provided to the decoder. This procedure was repeated for entire dataset  
782 ( $K=232$  fold) yielding both a binary a recognition score and the softmax probability  
783 assigned to each target class in the dataset. We define *classifier evidence* as the  
784 log-odds ratio of the probability assigned to the target category vs. the other categories  
785 in the dataset; this metric thus captures the degree to which a given drawing contains

786 visual information that is diagnostic of the target category (and not of the other  
787 categories in the dataset); these log-transformed values are also more suitable for the  
788 linear-mixed effects models used in analyses.

789 **CLIP Classifications.** For these analyses, we used the ViT-B/32  
790 implementation of CLIP publicly available at <https://github.com/openai/CLIP>. Model  
791 features were extracted for center-cropped versions of each sketch in the entire dataset  
792 ( $N=37770$ ), and for the tokenized, text versions of the labels for each of the 48  
793 categories (e.g. "a dog"). We then computed the cosine similarity between the features  
794 for each sketch and each of the 48 category labels and assessed which category label  
795 received the highest similarity. If the category label that had the highest similar was the  
796 category children were prompted to draw, this was counted as a correct classification.

797 **Human recognition scores: Experimental Dataset.** We measured the  
798 recognizability of each drawing in the controlled, experimental dataset via an online  
799 recognition experiment. Adult participants based in the U.S. were recruited via Prolific  
800 for a 15-minute experiment and asked to identify the category depicted in a random  
801 subset of approximately 140 drawings; each drawing was shown to 10 participants.  
802 Participants were shown these drawings in a random sequence and asked "What does  
803 this look like?" and selected their responses from the set of 12 categories and were  
804 encouraged to provide their best guess if they were unsure. No participants were  
805 excluded from analysis for missing the catch trial, which was included to verify that  
806 participants could accurately describe their goal in this task. We then computed a  
807 recognition score for each drawing, reflecting the proportion of participants who  
808 correctly identified the target category.

809 **Mixed-effect models.** Two mixed effects models were fit to assess the degree  
810 to which children produced more recognizable drawings across childhood. A first  
811 generalized mixed effect model was fit to the binary classification scores for each  
812 drawing using a logit linking function. A second linear mixed effect model was fit to the  
813 log-odds target probability assigned to each drawing, restricting our analyses to  
814 correctly classified drawings. In both cases, we included fixed effects of children's age

815 (in years), estimated drawing frequency for each category (via parental report), their  
816 interaction, children's estimated tracing score (see above), the time children spent  
817 drawing (in seconds), the mean intensity of the drawing (i.e. percentage of non-white  
818 pixels), and the number of strokes children used. All predictors were scaled to have a  
819 mean of 0 and a standard deviation of 1. Random intercepts were included for each  
820 participant and each category.

821 **Animacy & object size information in misclassified drawings.** For each  
822 misclassified drawing, we calculated whether the category assigned by the logistic  
823 regression classifier was the same animacy as the target category, assigning a binary  
824 animacy classification score for each drawing. The same procedure was repeated for  
825 inanimate objects with respect to their real-world object size (big objects: larger than a  
826 chair, small objects: can be held with one hand) (Konkle & Oliva, 2011; Long et al.,  
827 2016). These binary scores were averaged for each age and category, yielding a value  
828 between 0 and 1 representing the proportion of the drawings that were identified as  
829 having the correct animacy/size. As the proportion of animals/inanimate objects and  
830 big/small inanimate objects was not exactly balanced in the dataset, we subtracted the  
831 baseline prevalence for each broad category (i.e for animals, objects, big objects, and  
832 small objects) from this proportion. These values are plotted in Figures 3B,C, as are  
833 the bootstrapped 95% confidence intervals calculated using the baseline-corrected  
834 category values.

835 **Visual recognition task**

836 **Behavioral task.** On each trial of the guessing game, a photograph or drawing  
837 of an object category was presented on the screen, and children were asked to "tap the  
838 [animal/vehicle/object] that goes with the [drawing/picture]"; response choices were  
839 indicated by circular buttons that were filled photographs of canonical exemplars from  
840 each category, as well as the name of the category written above; the position of these  
841 response buttons was randomized for each participant. A fifth response choice was a  
842 button with a question-mark icon that could be used by participants to indicate they

843 didn't know which category the drawing belonged to. To familiarize participants with  
844 the interface, the first four trials of every game were four photograph trials, one for each  
845 of the response choices. To encourage accurate guessing, a pleasant sound was played  
846 when the correct category was chosen, and the box surrounding the image briefly  
847 turned green; no feedback was given for incorrect trials. Every ten trials, a catch trial  
848 appeared where participants were required to match a very similar photograph to the  
849 photographic response buttons.

850 **Drawing selection.** We selected four subsets of categories for the guessing  
851 game at the station: small animals (*dog, fish, rabbit, bird*), vehicles (*train, car, airplane,*  
852 *boat*), small, inanimate objects (*hat, bottle, cup, lamp*), and large animals (*camel, sheep,*  
853 *bear, tiger*). Each version of the guessing game ran separately for approximately two  
854 months. For each game, we randomly selected drawings (20-25 per category, depending  
855 on availability) made by children ages 4-8 at the drawing station. We chose this age  
856 range to cover a wide range of drawing abilities and to ensure equal numbers of  
857 drawings were included per age group (as 9-10 year-old's are infrequent visitors to the  
858 museum). This resulted in 516—616 drawings for each guessing game from which 48  
859 drawings were randomly sampled for each participant (8 drawings made by 4-, 5-, 6-, 7-,  
860 and 8-year-olds). If children completed the entire session, this resulted in a total of 48  
861 trials for each participant (40 drawing trials and 8 photograph matching trials).

862 **Recognition data inclusion.** As with the drawing data, we excluded any  
863 sessions where there was reported interference from parents or other children. As  
864 2-year-old's showed significantly better performance than 3-year-old's in our first two  
865 guessing games – signaling some interference from their caregivers or siblings that was  
866 not reported in the surveys – we chose to exclude 2-year-old's from subsequent analyses.  
867 We excluded children who started the game but did not complete more than 1 trial  
868 after the practice trials ( $N = 1068$  participants) and the 238 adults who participated.  
869 We also excluded all trials with reaction times slower than 10s or faster than 100ms,  
870 judging these to be off-task responses. Next, we excluded participants on the basis of  
871 their performance on practice and catch photograph matching trials. Given that these

872 catch trials presented a very easy recognition task, we excluded participants who did  
873 not achieve at least 75% accuracy on these trials ( $N = 795$ ). The remaining 1789  
874 participants who met this criterion completed an average of  $M=21.69$  trials. On total,  
875 we analyzed 36,615 trials where children recognized each other's drawings. These  
876 analysis choices were pre-registered after examining data from two of the guessing games  
877 and then applied to the entire dataset (see registrations on <https://osf.io/qymjr/>).

878 **Recognition data analyses.** To calculate the classifier evidence associated  
879 with each sketch that children recognized, we used the same visual encoder to extract  
880 visual features for each sketch (see *Visual Encoder*), and iteratively trained logistic  
881 regression classifiers (see *Logistic Regression Classifier*). For these analyses, we  
882 restricted the classification set to the drawings that were presented in each version of  
883 the guessing game to match the task conditions of the guessing game. We trained a  
884 separate logistic regression for each sketch that was presented using leave-one-out  
885 cross-validation. This procedure thus yielded probabilities assigned to each of four  
886 categories in each guessing game; these probabilities were used to calculate the log-odds  
887 ratios for the target category of each sketch which we refer to as *classifier evidence*. Due  
888 to random sampling, not every sketch included in the game had valid guesses associated  
889 with it; these sketches were thus not included in analyses. We then modeled children's  
890 recognition behavior in a generalized linear mixed-effect model, where recognizer age (in  
891 years), classifier evidence, and their interaction were specified as fixed effects. All  
892 predictors were scaled between 0 and 1. We included random intercepts for the intended  
893 category of the sketch and for each subject who participated in the guessing game;  
894 random slopes were also included for the effect of classifier evidence on each intended  
895 category.

896 **Semantic part annotation task**

897 **Crowdsourcing category part decompositions.** We designed a web-based  
898 crowdsourcing platform and recruited 50 English-speaking adult participants from  
899 Prolific to identify the basic parts of objects for each of the 16 object categories. On

900 each trial, participants were cued with a text label of an object category and asked to  
901 list 3 to 10 object parts that came to mind (e.g., head, leg, tail, etc. for “tiger”).  
902 Participants were instructed to write only concrete parts of an object (e.g., “tail”)  
903 rather than abstract attributes (e.g., “tufted”), to use common names of parts rather  
904 than technical jargon (e.g., “prehensile”), and to generate as complete a part list as  
905 they could for each object category. We applied lemmatization to the resulting part  
906 decompositions to remove redundant part labels, such as “hoof” and “hoofs”, and  
907 manually edited part labels that were spelled incorrectly or with alternative spellings.  
908 We then selected the top 10% of part names that were most frequently listed. This  
909 generated a total of 82 object parts with a range of 5-13 possible parts per object  
910 category.

911 **Part labeling task.** We developed a web-based annotation paradigm adapted  
912 from previous research (Huey, Walker, & Fan, 2021; Mukherjee, Hawkins, & Fan, 2019)  
913 to obtain detailed annotations of how each pen stroke in children’s drawings  
914 corresponded to the different parts of the depicted objects. 1,034 English-speaking  
915 adult participants were recruited from Prolific and completed the semantic annotation  
916 task. We excluded data from 78 additional participants for experiencing technical  
917 difficulties with the web interface ( $N=11$ ) and for having low accuracy on our  
918 attention-check trial ( $N=67$ ). Data collection was stopped when every drawing had  
919 received annotations from at least three annotators.

920 Each annotator was presented with a set of 8 drawings randomly sampled from  
921 the drawing dataset but consistent within the same animacy and object size (i.e., small  
922 animals, large animals, vehicles, household objects). Each drawing was accompanied by  
923 the name of its object category (e.g., “airplane”), as well as a gallery of crowd-sourced  
924 part labels that corresponded to it. For each stroke in the presented drawing,  
925 annotators were prompted to tag it with the part label that described the part of the  
926 depicted object that it represented. Annotators were permitted to label a stroke with  
927 multiple part labels if they believed a stroke to represent multiple different parts of the  
928 depicted object, and were able to write their own custom label if they believed that

929 none of the provided part labels were fitting. They could also label a stroke as  
930 unintelligible if they could not discern what it represented. Annotators also completed  
931 an “attention-check” trial, consisting of a pre-selected drawing that had been annotated  
932 by a researcher and then randomly inserted into the set of drawings. If annotators did  
933 not match the researcher’s annotation criteria for this drawing, data sessions from these  
934 annotators were excluded from subsequent analysis.

935 **Annotation data preprocessing.** First, we evaluated how often annotators  
936 agreed on what each stroke of children’s drawings represented by calculating the  
937 inter-rater consistency among annotators. Across drawings, annotators agreed on the  
938 same part label for 69.9% of strokes. There was modest improvement in agreement  
939 across age, with drawings produced by older children eliciting more consistent  
940 annotations (4-year-old drawings = 68.3% mean agreement, 8-year-old drawings =  
941 69.8% mean agreement). We retained stroke annotations that were assigned the same  
942 part label(s) by at least two of three annotators. While annotators infrequently wrote  
943 custom labels (we did not analyze custom annotations for the present analysis) they  
944 only used 68 of the available 82 part labels. Our resultant dataset therefore contained  
945 14,159 annotated strokes across 2,088 drawings.

946 **Part inclusion and emphasis calculation.** For part inclusion, we calculated  
947 the number of unique object parts assigned to each drawing; strokes labeled as  
948 unintelligible were not counted as distinct parts. For part emphasis, we calculated the  
949 proportion of the total length of strokes that were attributed to a particular object part  
950 in a drawing (e.g., wings), relative the total length of all strokes in the entire drawing  
951 (including strokes that were not agreed upon or that were unintelligible). If strokes  
952 were used to represent multiple object parts, we took the total length of the stroke and  
953 divided it by the number of parts that it was assigned to.

## 954 **Data Availability**

955 Source data are provided with this paper. The drawings and pre-processed data  
956 that support the findings are available at <https://osf.io/qymjr/>

**957 Code Availability**

958 The code used analyze the data are available at <https://osf.io/qymjr/>

**959 Acknowledgements**

960 We gratefully acknowledge the San Jose Children's Discovery Museum for their  
961 collaboration and for hosting the drawing station where these data were collected. We  
962 also thank the members of the Stanford Language and Cognition lab for their feedback  
963 on this project throughout several years and Megan Merrick for assistance with data  
964 collection. This work was funded by an NSF SPRF-FR Grant #1714726 to BLL, a NIH  
965 K99 HD108386 to BLL, and a Jacobs Foundation Fellowship to MCF.

966

## References

- 967 Bainbridge, W. A., Hall, E. H., & Baker, C. I. (2019). Drawings of real-world scenes  
968 during free recall reveal detailed object and spatial information in memory.  
969 *Nature communications*, 10(1), 1–13.
- 970 Balas, B., & Saville, A. (2020). Neural sensitivity to natural image statistics changes  
971 during middle childhood. *Developmental Psychobiology*.
- 972 Barrett, M., & Light, P. (1976). Symbolism and intellectual realism in children's  
973 drawings. *British Journal of Educational Psychology*, 46(2), 198–202.
- 974 Battleday, R. M., Peterson, J. C., & Griffiths, T. L. (2020). Capturing human  
975 categorization of natural images by combining deep networks and cognitive  
976 models. *Nature communications*, 11(1), 1–14.
- 977 Biederman, I., & Ju, G. (1988). Surface versus edge-based determinants of visual  
978 recognition. *Cognitive psychology*, 20(1), 38–64.
- 979 BLINDED. (in press). Developmental changes in drawing production under different  
980 memory demands in a u.s. and chinese sample. *Developmental Psychology*.  
981 Retrieved from BLINDED
- 982 Bova, S. M., Fazzi, E., Giovenzana, A., Montomoli, C., Signorini, S. G., Zoppello, M., &  
983 Lanzi, G. (2007). The development of visual object recognition in school-age  
984 children. *Developmental neuropsychology*, 31(1), 79–102.
- 985 Bozeat, S., Lambon Ralph, M. A., Graham, K. S., Patterson, K., Wilkin, H., Rowland,  
986 J., . . . Hodges, J. R. (2003). A duck with four legs: Investigating the structure of  
987 conceptual knowledge using picture drawing in semantic dementia. *Cognitive  
988 neuropsychology*, 20(1), 27–47.
- 989 Braginsky, M., Yurovsky, D., Marchman, V. A., & Frank, M. C. (2019). Consistency  
990 and variability in children's word learning across languages. *Open Mind*, 3, 52–67.
- 991 Bremner, J. G., & Moore, S. (1984). Prior visual inspection and object naming: Two  
992 factors that enhance hidden feature inclusion in young children's drawings. *British  
993 Journal of Developmental Psychology*, 2(4), 371–376.
- 994 Carey, S., & Bartlett, E. (1978). Acquiring a single new word. *ERIC*.

- 995 Cavanagh, P. (2005). The artist as neuroscientist. *Nature*, 434(7031), 301–307.
- 996 Chamberlain, R., Drake, J. E., Kozbelt, A., Hickman, R., Siev, J., & Wagemans, J.  
997 (2019). Artists as experts in visual cognition: An update. *Psychology of  
998 Aesthetics, Creativity, and the Arts*, 13(1), 58.
- 999 Chamberlain, R., Kozbelt, A., Drake, J. E., & Wagemans, J. (2021). Learning to see by  
1000 learning to draw: A longitudinal analysis of the relationship between  
1001 representational drawing training and visuospatial skill. *Psychology of Aesthetics,  
1002 Creativity, and the Arts*, 15(1), 76.
- 1003 Clottes, J. (2008). *Cave art*. Phaidon London.
- 1004 Cohen, M. A., Dilks, D. D., Koldewyn, K., Weigelt, S., Feather, J., Kell, A. J., ...  
1005 others (2019). Representational similarity precedes category selectivity in the  
1006 developing ventral visual pathway. *NeuroImage*, 197, 565–574.
- 1007 Conwell, C., Prince, J. S., Kay, K. N., Alvarez, G. A., & Konkle, T. (2022). What can  
1008 1.8 billion regressions tell us about the pressures shaping high-level visual  
1009 representation in brains and machines? *bioRxiv*, 2022–03.
- 1010 Cox, M. V., & Ralph, M. L. (1996). Young children's ability to adapt their drawings of  
1011 the human figure. *Educational Psychology*, 16(3), 245–255.
- 1012 Dekker, T., Mareschal, D., Sereno, M. I., & Johnson, M. H. (2011). Dorsal and ventral  
1013 stream activation and object recognition performance in school-age children.  
1014 *NeuroImage*, 57(3), 659–670.
- 1015 DeLoache, J. S., Pierroutsakos, S. L., & Uttal, D. H. (2003). The origins of pictorial  
1016 competence. *Current Directions in Psychological Science*, 12(4), 114–118.
- 1017 Deregowksi, J. B. (1989). Real space and represented space: Cross-cultural  
1018 perspectives. *Behavioral and Brain Sciences*, 12(1), 51–74.
- 1019 Fan, J. E., Wammes, J. D., Gunn, J. B., Yamins, D. L., Norman, K. A., &  
1020 Turk-Browne, N. B. (2020). Relating visual production and recognition of objects  
1021 in human visual cortex. *Journal of Neuroscience*, 40(8), 1710–1721.
- 1022 Fan, J. E., Yamins, D. L., & Turk-Browne, N. B. (2018). Common object  
1023 representations for visual production and recognition. *Cognitive science*, 42(8),

- 1024 2670–2698.
- 1025 Finke, R. A., & Slayton, K. (1988). Explorations of creative visual synthesis in mental  
1026 imagery. *Memory & cognition*, 16(3), 252–257.
- 1027 Fisher, A. V., Godwin, K. E., & Matlen, B. J. (2015). Development of inductive  
1028 generalization with familiar categories. *Psychonomic bulletin & review*, 22,  
1029 1149–1173.
- 1030 Foster-Hanson, E., & Rhodes, M. (2019). Is the most representative skunk the average  
1031 or the stinkiest? developmental changes in representations of biological categories.  
1032 *Cognitive psychology*, 110, 1–15.
- 1033 Freeman, N. H. (1987). Current problems in the development of representational  
1034 picture-production. *Archives de psychologie*.
- 1035 Freeman, N. H., & Janikoun, R. (1972). Intellectual realism in children's drawings of a  
1036 familiar object with distinctive features. *Child development*, 1116–1121.
- 1037 Fury, G., Carlson, E. A., & Sroufe, A. (1997). Children's representations of attachment  
1038 relationships in family drawings. *Child development*, 68(6), 1154–1164.
- 1039 Gibson, J. J. (1971). The information available in pictures. *Leonardo*, 27–35.
- 1040 Gomez, J., Natu, V., Jeska, B., Barnett, M., & Grill-Spector, K. (2018). Development  
1041 differentially sculpts receptive fields across early and high-level human visual  
1042 cortex. *Nature communications*, 9(1), 788.
- 1043 Goodenough, F. L. (1963). *Goodenough-harris drawing test*. Harcourt Brace Jovanovich  
1044 New York.
- 1045 Goodman, J. C., Dale, P. S., & Li, P. (2008). Does frequency count? parental input  
1046 and the acquisition of vocabulary. *Journal of Child Language*, 35(3), 515-531.
- 1047 Gregory, R. L. (1973). *Eye and brain: The psychology of seeing*. McGraw-Hill.
- 1048 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). Most people are not weird. *Nature*,  
1049 466(7302), 29–29.
- 1050 Hertzmann, A. (2020). Why do line drawings work? a realism hypothesis. *Perception*,  
1051 49(4), 439–451.
- 1052 Huey, H., Walker, C. M., & Fan, J. E. (2021). How do the semantic properties of visual

- 1053 explanations guide causal inference? In *Proceedings of the annual meeting of the*  
1054 *cognitive science society* (Vol. 43).
- 1055 Itakura, S. (1994). Recognition of line-drawing representations by a chimpanzee (pan  
1056 troglodytes). *The Journal of general psychology*, 121(3), 189–197.
- 1057 James, K. H. (2017). The importance of handwriting experience on the development of  
1058 the literate brain. *Current Directions in Psychological Science*, 26(6), 502–508.
- 1059 Juttner, M., Muller, A., & Rentschler, I. (2006). A developmental dissociation of  
1060 view-dependent and view-invariant object recognition in adolescence. *Behavioural*  
1061 *brain research*, 175(2), 420–424.
- 1062 Juttner, M., Wakui, E., Petters, D., & Davidoff, J. (2016). Developmental  
1063 commonalities between object and face recognition in adolescence. *Frontiers in*  
1064 *psychology*, 7.
- 1065 Karmiloff-Smith, A. (1990). Constraints on representational change: Evidence from  
1066 children's drawing. *Cognition*, 34(1), 57–83.
- 1067 Kellogg, R. (1969). *Analyzing children's art*. National Press Books Palo Alto, CA.
- 1068 Kennedy, J. M., & Ross, A. S. (1975). Outline picture perception by the songe of  
1069 papua. *Perception*, 4(4), 391–406.
- 1070 Kersey, A. J., Clark, T. S., Lussier, C. A., Mahon, B. Z., & Cantlon, J. F. (2015).  
1071 Development of tool representations in the dorsal and ventral visual object  
1072 processing pathways. *Cerebral Cortex*, 26(7), 3135–3145.
- 1073 Konkle, T., & Oliva, A. (2011). Canonical visual size for real-world objects. *Journal of*  
1074 *experimental psychology: human perception and performance*, 37(1), 23.
- 1075 Kozbelt, A. (2001). Artists as experts in visual cognition. *Visual cognition*, 8(6),  
1076 705–723.
- 1077 Lewis, M., Balamurugan, A., Zheng, B., & Lupyan, G. (2021). Characterizing  
1078 variability in shared meaning through millions of sketches. In *Proceedings of the*  
1079 *annual meeting of the cognitive science society* (Vol. 43).
- 1080 Li, J. X., & James, K. H. (2016). Handwriting generates variable visual output to  
1081 facilitate symbol learning. *Journal of Experimental Psychology: General*, 145(3),

- 1082 298.
- 1083 Long, B., Fan, J., & Frank, M. C. (2018). Drawings as a window into developmental  
1084 changes in object representations. In *Proceedings of the 40th annual meeting of*  
1085 *the cognitive science society*.
- 1086 Long, B., Konkle, T., Cohen, M. A., & Alvarez, G. A. (2016). Mid-level perceptual  
1087 features distinguish objects of different real-world sizes. *Journal of Experimental*  
1088 *Psychology: General*, 145(1), 95.
- 1089 Long, B., Moher, M., Carey, S. E., & Konkle, T. (2019). Animacy and object size are  
1090 reflected in perceptual similarity computations by the preschool years. *Visual*  
1091 *Cognition*, 27(5-8), 435–451.
- 1092 Long, B., Störmer, V. S., & Alvarez, G. A. (2017). Mid-level perceptual features  
1093 contain early cues to animacy. *Journal of Vision*, 17(6), 20–20.
- 1094 Long, B., Yu, C.-P., & Konkle, T. (2018). Mid-level visual features underlie the  
1095 high-level categorical organization of the ventral stream. *Proceedings of the*  
1096 *National Academy of Sciences*, 115(38), E9015–E9024.
- 1097 Longcamp, M., Zerbato-Poudou, M.-T., & Velay, J.-L. (2005). The influence of writing  
1098 practice on letter recognition in preschool children: A comparison between  
1099 handwriting and typing. *Acta psychologica*, 119(1), 67–79.
- 1100 Luquet, G.-H. (1927). Le dessin enfantin.(bibliothèque de psychologie de l" enfant et de  
1101 pédagogie.).
- 1102 Mash, C. (2006). Multidimensional shape similarity in the development of visual object  
1103 classification. *Journal of Experimental Child Psychology*, 95(2), 128–152.
- 1104 Minsky, M., & Papert, S. (1972). *Artificial intelligence progress report* (Tech. Rep.).  
1105 Cambridge, MA, USA.
- 1106 Morra, S., & Panesi, S. (2017). From scribbling to drawing: the role of working  
1107 memory. *Cognitive Development*, 43, 142–158.
- 1108 Mukherjee, K., Hawkins, R. X., & Fan, J. W. (2019). Communicating semantic part  
1109 information in drawings. In *Cogsci* (pp. 2413–2419).
- 1110 Natu, V. S., Barnett, M. A., Hartley, J., Gomez, J., Stigliani, A., & Grill-Spector, K.

- 1111       (2016). Development of neural sensitivity to face identity correlates with  
1112       perceptual discriminability. *Journal of Neuroscience*, 36(42), 10893–10907.
- 1113   Nishimura, M., Scherf, K. S., Zachariou, V., Tarr, M. J., & Behrmann, M. (2015). Size  
1114       precedes view: developmental emergence of invariant object representations in  
1115       lateral occipital complex. *Journal of cognitive neuroscience*, 27(3), 474–491.
- 1116   Nishimura, M., Scherf, S., & Behrmann, M. (2009). Development of object recognition  
1117       in humans. *F1000 biology reports*, 1.
- 1118   Perdreau, F., & Cavanagh, P. (2013a). The artist's advantage: Better integration of  
1119       object information across eye movements. *i-Perception*, 4(6), 380–395.
- 1120   Perdreau, F., & Cavanagh, P. (2013b). Is artists' perception more veridical? *Frontiers*  
1121       in *neuroscience*, 7, 6.
- 1122   Perdreau, F., & Cavanagh, P. (2014). Drawing skill is related to the efficiency of  
1123       encoding object structure. *i-Perception*, 5(2), 101–119.
- 1124   Pereira, A. F., & Smith, L. B. (2009). Developmental changes in visual object  
1125       recognition between 18 and 24 months of age. *Developmental science*, 12(1),  
1126       67–80.
- 1127   Piaget, J. (1929). The child's concept of the world. *Londres, Routledge & Kegan Paul*.
- 1128   Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., . . . others  
1129       (2021). Learning transferable visual models from natural language supervision. In  
1130       *International conference on machine learning* (pp. 8748–8763).
- 1131   Rehrig, G., & Stromswold, K. (2018). What does the dap: Iq measure?: Drawing  
1132       comparisons between drawing performance and developmental assessments. *The*  
1133       *Journal of genetic psychology*, 179(1), 9–18.
- 1134   Rosch, E. (1978). Principles of categorization.
- 1135   Sandkühler, R., Jud, C., Andermatt, S., & Cattin, P. C. (2018). Airlab: Autograd  
1136       image registration laboratory. *arXiv preprint arXiv:1806.09907*.
- 1137   Sayim, B. (2011, October). What line drawings reveal about the visual brain. , 1–4.
- 1138   Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale  
1139       image recognition. *arXiv preprint arXiv:1409.1556*.

- 1140 Sitton, R., & Light, P. (1992). Drawing to differentiate: Flexibility in young children's  
1141 human figure drawings. *British Journal of Developmental Psychology*, 10(1),  
1142 25–33.
- 1143 Soja, N. N., Carey, S., & Spelke, E. S. (1991). Ontological categories guide young  
1144 children's inductions of word meaning: Object terms and substance terms.  
1145 *Cognition*, 38(2), 179–211.
- 1146 Tanaka, M. (2007). Recognition of pictorial representations by chimpanzees (pan  
1147 troglodytes). *Animal cognition*, 10(2), 169–179.
- 1148 Tchalenko, J. (2009). Segmentation and accuracy in copying and drawing: Experts and  
1149 beginners. *Vision research*, 49(8), 791–800.
- 1150 Tversky, B. (1985). Development of taxonomic organization of named and pictured  
1151 categories. *Developmental Psychology*, 21(6), 1111.
- 1152 Tversky, B. (1989). Parts, partonomies, and taxonomies. *Developmental Psychology*,  
1153 25(6), 983.
- 1154 Tversky, B., & Hemenway, K. (1984). Objects, parts, and categories. *Journal of  
1155 experimental psychology: General*, 113(2), 169.
- 1156 Vales, C., Stevens, P., & Fisher, A. V. (2020). Lumping and splitting: Developmental  
1157 changes in the structure of children's semantic networks. *Journal of Experimental  
1158 Child Psychology*, 199, 104914.
- 1159 Weigelt, S., Koldewyn, K., Dilks, D. D., Balas, B., McKone, E., & Kanwisher, N.  
1160 (2014). Domain-specific development of face memory but not face perception.  
1161 *Developmental Science*, 17(1), 47–58.
- 1162 Yamins, D., Hong, H., Cadieu, C., Solomon, E., Seibert, D., & DiCarlo, J. (2014).  
1163 Performance-optimized hierarchical models predict neural responses in higher  
1164 visual cortex. *Proceedings of the National Academy of Sciences*, 111(23),  
1165 8619–8624.

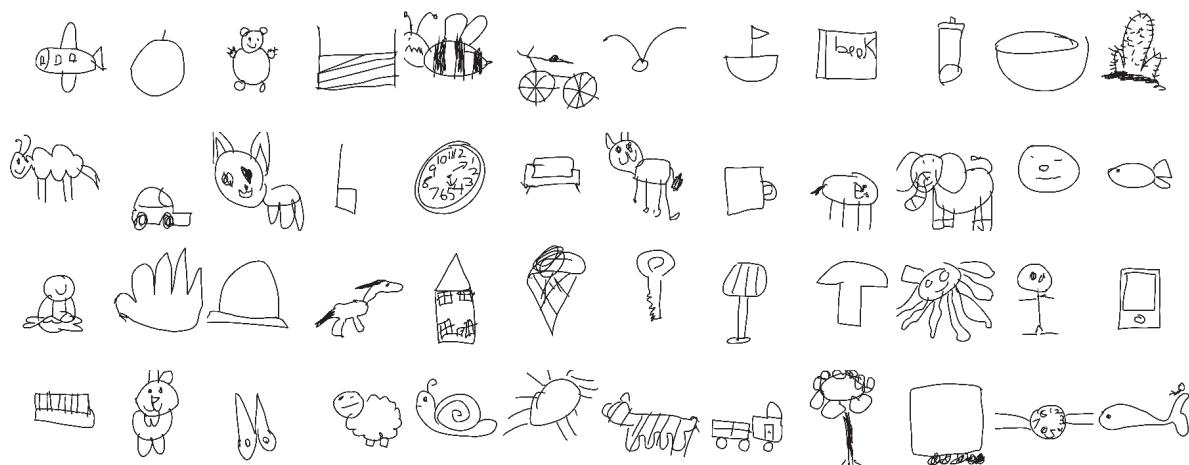
## Appendix A

Supplemental analyses: visual production across childhood

Age	Number of participants	Number of drawings
2-year-olds	1231	3651
3-year-olds	1402	5342
4-year-olds	1451	6559
5-year-olds	1189	6411
6-year-olds	878	4990
7-year-olds	660	3817
8-year-olds	478	2570
9-year-olds	309	1800
10+-year-olds	486	2630

Table A1

*Number of participants and drawings included in the filtered drawing dataset by each age group.*

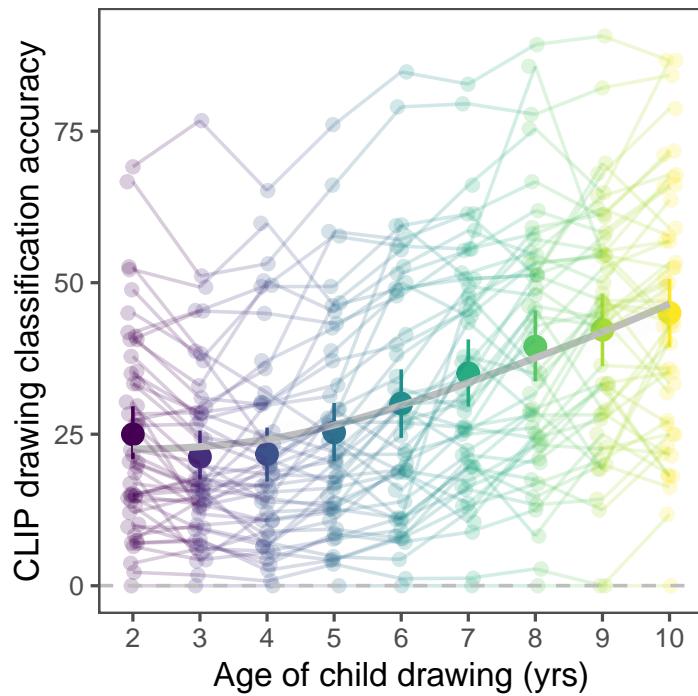


*Figure A1.* Examples of correctly classified drawings from each of the 48 categories presented at the experiment station in alphabetical order: airplane, apple, bear, bed, bee, bike, bird, boat, book, bottle, bowl, cactus, (2nd row): camel, car, cat, chair, clock, couch, cow, cup, dog, elephant, face, fish, (3rd row): frog, hand, hat, horse, house, ice cream, key, lamp, mushroom, octopus, person, phone, (4th row): piano, rabbit, scissors, sheep, snail, spider, tiger, train, tree, TV, watch.

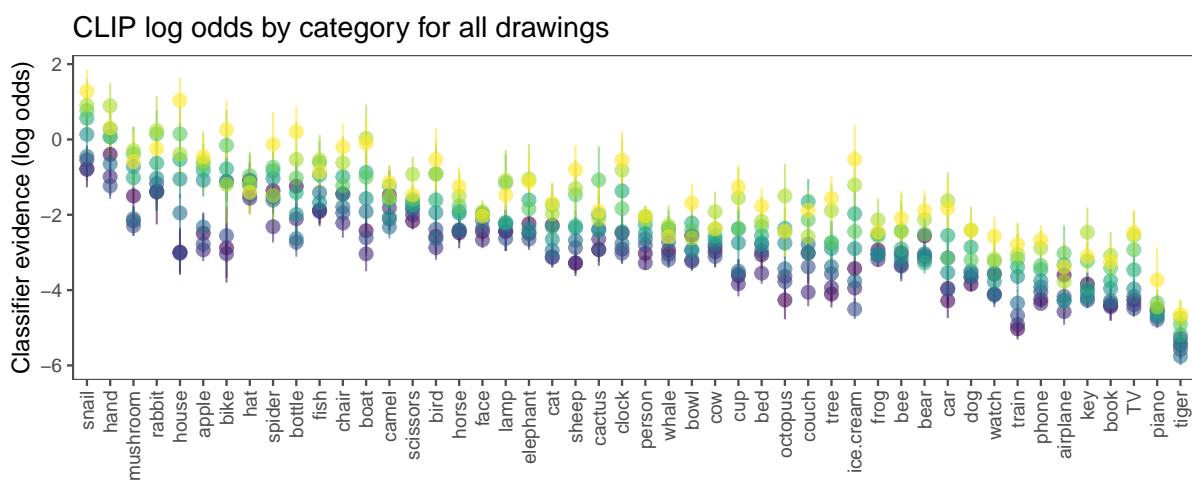
	Term	VIF	SE_factor
1	Age	1.26	1.12
2	Drawing Frequency	1.00	1.00
3	Tracing score	1.24	1.11
4	Draw duration	1.30	1.14
5	Ink used	1.29	1.14
6	Number of strokes	1.07	1.03
7	Age* Drawing frequency	1.01	1.00

Table A2

*Results of the multicollinearity analysis for the predictors used in the main GLMM predicting the recognizability of children's drawings.*



*Figure A2.* Figure 2A, redone with CLIP classifications. Y-axis shows classification accuracy as a function of children's age (x-axis). Each dot represents data from an individual category, which are connected by individual trend lines. Error bars represent bootstrapped 95% confidence intervals.

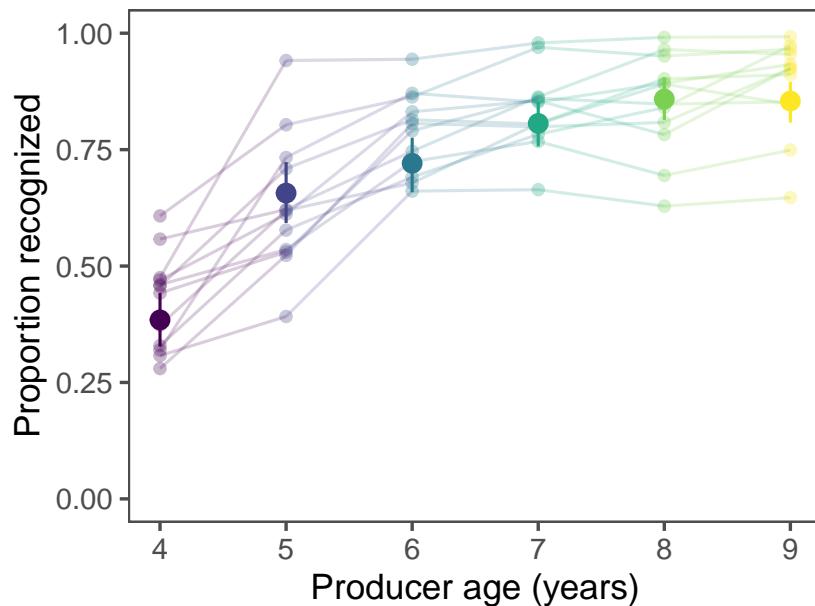


*Figure A3.* CLIP log-odds probabilities (y-axis) assigned to each category as a function of children's age; each dot represents data from an individual category and age.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.319	0.178	-7.410	<0.001
Age (in years)	0.329	0.019	17.225	<0.001
Est. drawing frequency	0.274	0.177	1.551	0.121
Tracing score	0.279	0.020	14.320	<0.001
Time spent drawing	0.195	0.019	10.065	<0.001
'Ink' used	0.047	0.018	2.642	0.007
Number of strokes	0.070	0.030	2.338	0.019
Age*drawing frequency	0.029	0.014	2.030	0.042

Table A3

*Model coefficients of a GLMM predicting the recognizability of each drawing (i.e. binary classification scores) from CLIP classifications, including random intercepts for each category and participant.*

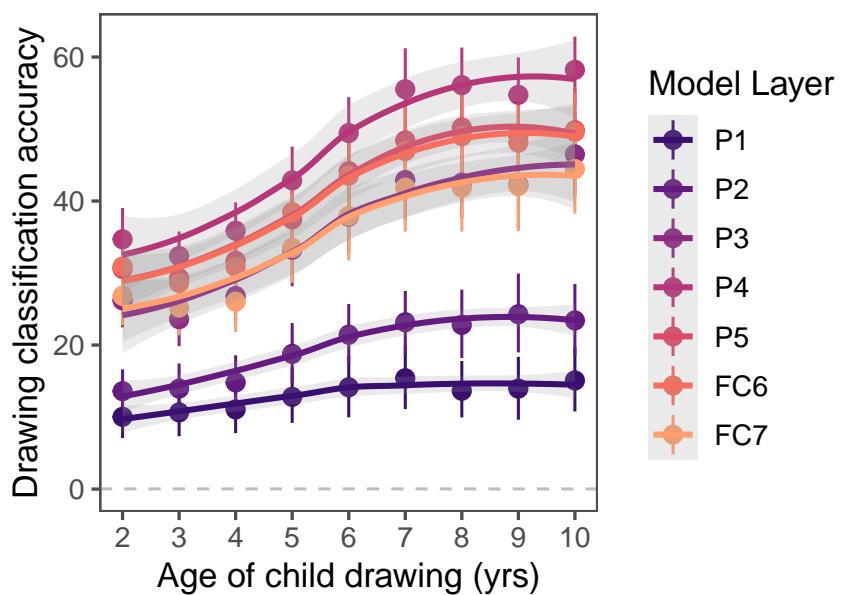


*Figure A4.* Drawing accuracy as a function of children's age; children drew in response to verbal prompts in a controlled, experimental setting. Y-axis reflects the proportion of human observers who correctly identified the drawing in a 12AFC guessing task. Error bars reflect 95 percent bootstrapped confidence intervals.

	Estimate	Std. Error	df	t value	Pr(> t )
(Intercept)	0.736	0.024	12.391	30.226	<0.001
Drawing frequency	0.043	0.024	9.853	1.833	0.097
Age (in years)	0.151	0.011	247.928	14.348	<0.001
Drawing frequency * Age	0.007	0.007	1282.435	1.073	0.283

Table A4

*Model coefficients of a linear mixed effect model predicting the recognizability of each drawing (as assessed by crowd-sourced adult behavioral data). Drawings were produced in an experimental context. All predictors are z-scored and random intercepts for each category and participant are included.*



*Figure A5.* Drawing accuracy as a function of children's age using embeddings from each layer in the VGG-19 network. Error bars represent 95 percent bootstrapped confidence intervals.

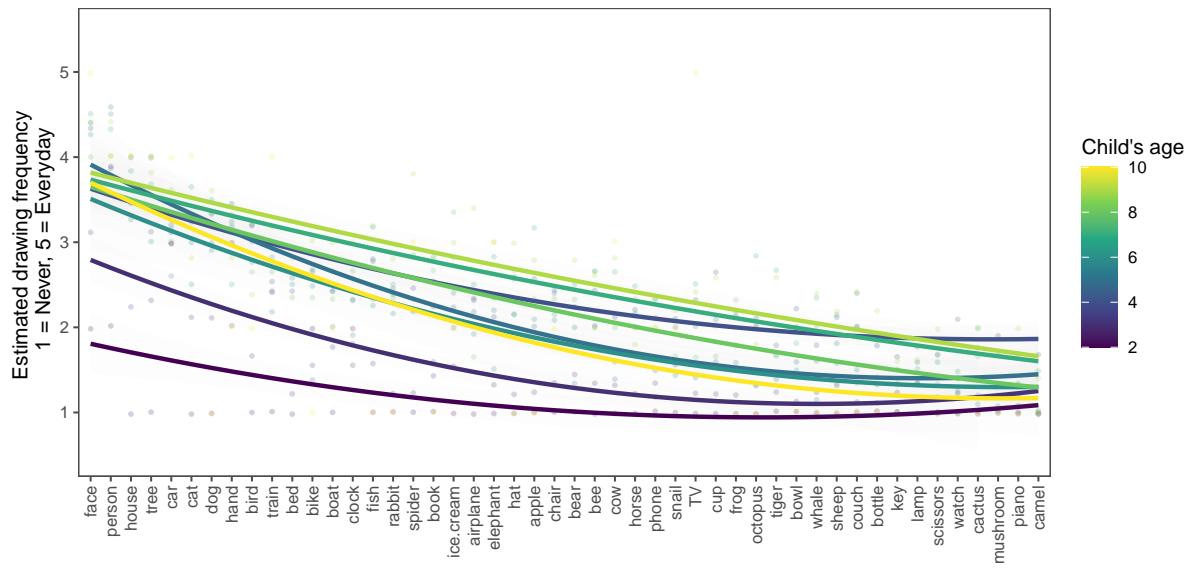


Figure A6. Frequency (y-axis) with which parents (recruited online) estimated their children drew each of the 48 categories in the dataset.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.965	0.237	-4.065	<0.001
Children's age	0.360	0.021	17.006	<0.001
Freq. in adult books	0.515	0.312	1.653	0.098
Est. AoA	0.702	0.418	1.678	0.093
Freq. in CHILDES	0.129	0.398	0.323	0.747
Drawing frequency	-0.290	0.326	-0.889	0.374

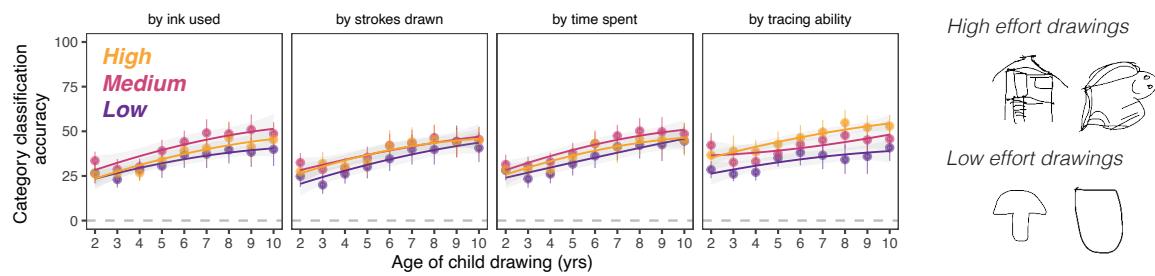
Table A5

*Model coefficients of a GLMM predicting the recognizability of each drawing (i.e. binary classification scores) from children's age, the frequency of each category in CHILDES, the estimated Age-of-Acquisition, and the frequency of each word in adult English books. All predictors are z-scored; random intercepts for each category and participant are included.*

	Estimate	Std. Error	df	t value	Pr(> t )
(Intercept)	-0.687	0.107	43.356	-6.434	<0.001
Children's age	0.111	0.015	3544.182	7.354	<0.001
Drawing frequency	0.020	0.114	42.904	0.174	0.863
Average tracing rating	0.101	0.015	3964.435	6.753	<0.001
Time spent drawing	-0.019	0.018	7710.462	-1.060	0.289
Ink spent	0.018	0.017	7888.865	1.042	0.297
Number of strokes	0.063	0.018	8110.693	3.535	0.000
Age*Drawing frequency	0.011	0.013	7672.028	0.811	0.417

Table A6

*Model coefficients of a linear mixed effect model predicting the log-odds probability assigned to correctly classified drawings using VGG embeddings, including random intercepts for each category and participant.*



*Figure A7.* (Left): Classification accuracy by age, split into bins according to whether children expended a greater/lesser amount of strokes, ink, or time, and by their estimated tracing abilities (see Methods). (Right): Example drawings where children spent higher/lower amounts of *effort*—greater/lower than average number of strokes, time spent drawing, or 'ink' used.

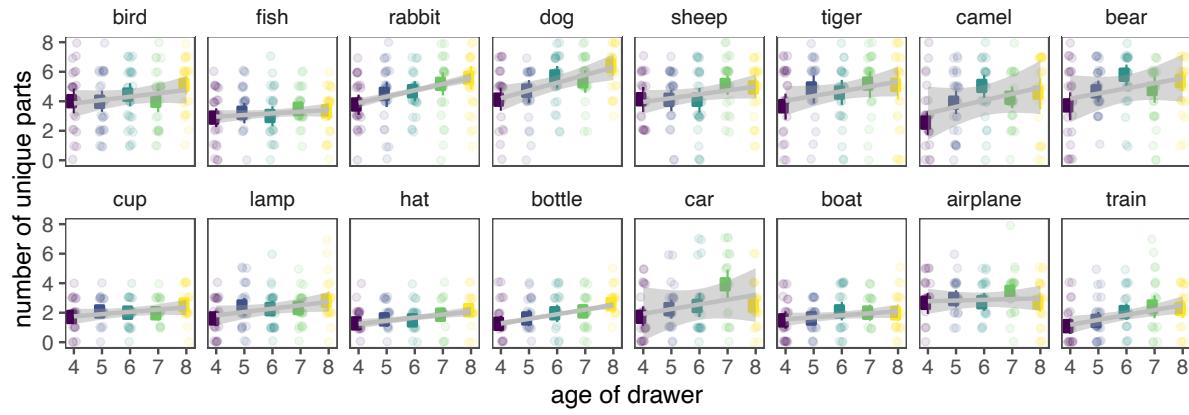


Figure A8. Number of unique parts per each object category included in the semantic part annotation subset ( $N=2,088$  drawings of 16 categories) across age.

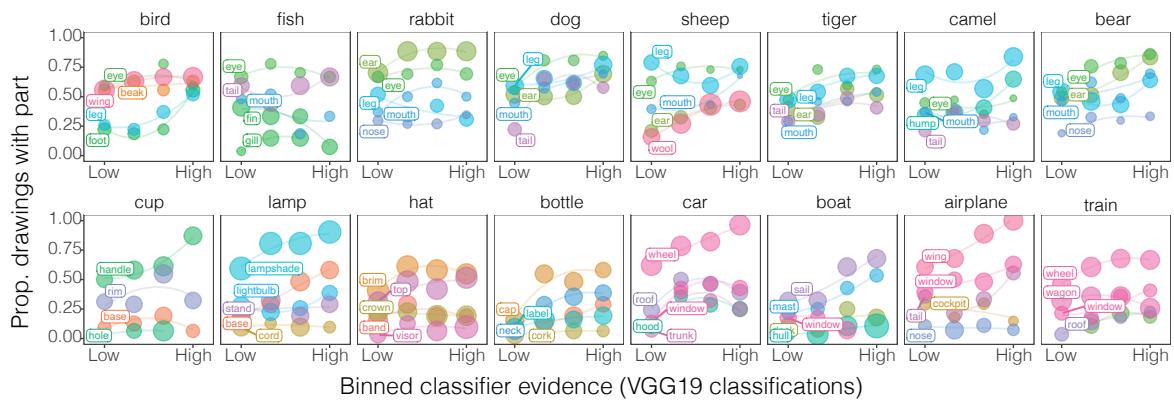
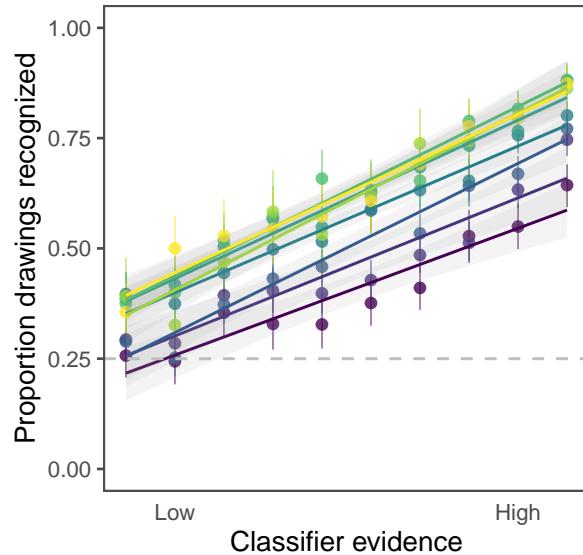


Figure A9. Object parts included for each category as a function of the VGG-19 classification evidence, binned into quartiles; as in the main texts, only the top 4 object parts that were frequently included (excluding head/body) are shown here. Dot size represents the visual emphasis on each part.

Appendix B  
Supplemental materials: visual recognition



*Figure B1.* Replication of the main interaction on visual recognition behaviors (proportion recognized, y-axis) by recognizer age (individual lines colored by age) and classifier evidence, here using CLIP classification probabilities (binned into deciles on the x-axis).

<sup>1166</sup> **Including only high-performing children.**

<sup>1167</sup> To ensure that these results were not driven by differences in motivation or  
<sup>1168</sup> general task performance, we also conducted our main analyses on a very restricted  
<sup>1169</sup> subset of our participants. We excluded any participant that did not achieve 100% on  
<sup>1170</sup> the photograph matching trials or that scored less than 50% on the drawing recognition  
<sup>1171</sup> trials. While this excluded nearly two-thirds of our participants, there were nonetheless  
<sup>1172</sup>  $N=649$  participants in this subset. Nonetheless, we still found the same pattern of  
<sup>1173</sup> results (see Table B1): older children were still better at recognizing drawings and at  
<sup>1174</sup> using diagnostic visual information in these drawings when recognizing them.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	0.668	0.099	6.715	<0.001
Classifier evidence (scaled)	0.518	0.051	10.057	<0.001
Recognizer age (scaled)	0.141	0.023	6.190	<0.001
Classifier evidence*Recognizer Age	0.056	0.023	2.464	0.014

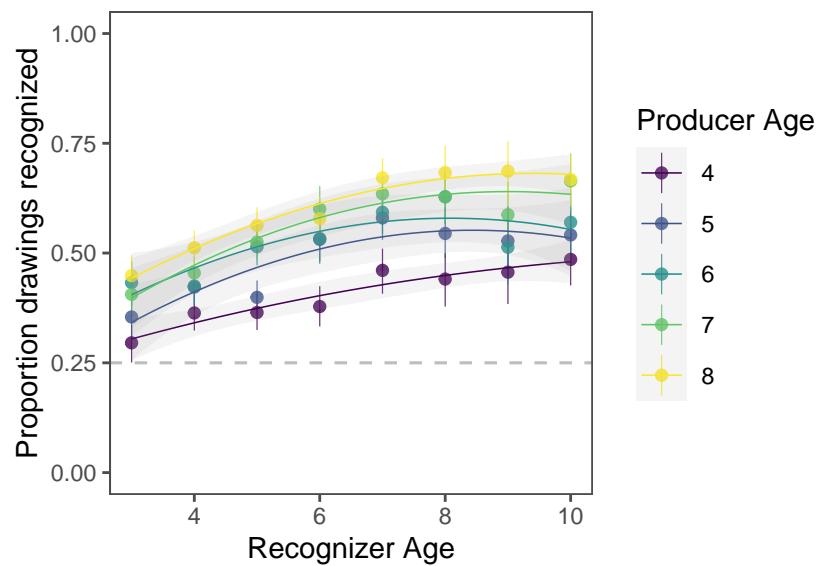
Table B1

*Model coefficients of a GLMM predicting visual recognition performance, excluding any participant who missed even one of the photograph trials or who scored less than 50% on drawing recognition trials.*

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	0.397	0.183	2.167	0.030
Classifier evidence (scaled)	0.987	0.215	4.590	<0.001
Recognizer age (scaled)	0.367	0.021	17.644	<0.001
Classifier evidence*Recognizer Age	0.091	0.018	5.079	<0.001

Table B2

*Model coefficients of a GLMM predicting visual recognition performance as a function of recognizer age and the category-diagnostic information in drawings (derived from CLIP embeddings, see Methods).*

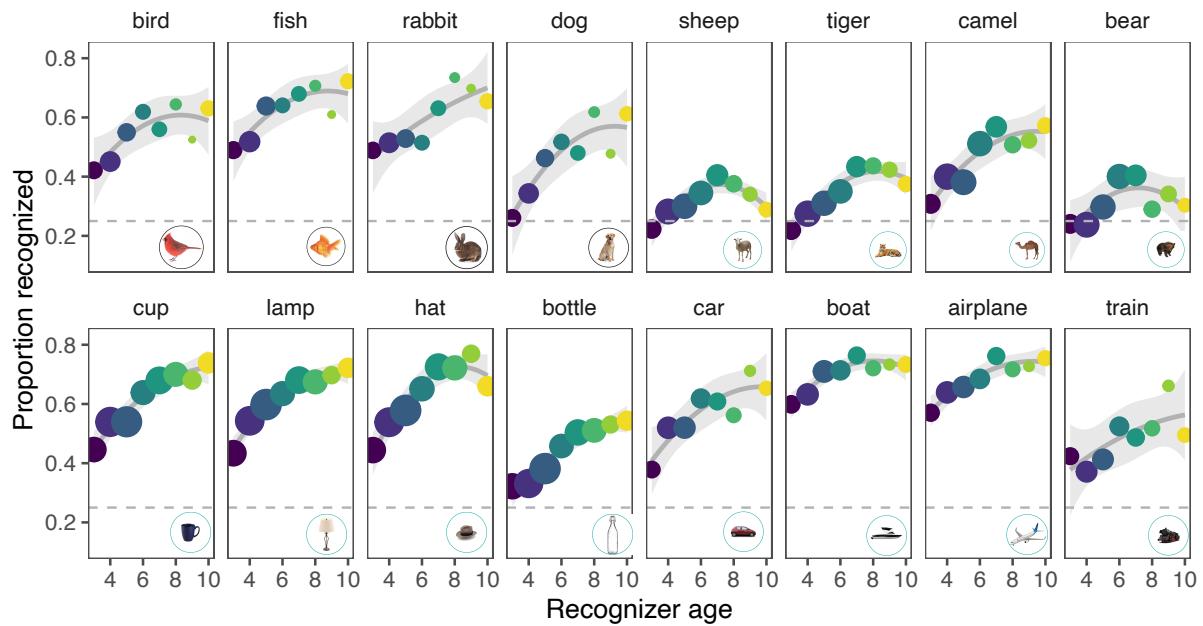


*Figure B2.* Proportion of drawings recognized (y-axis) as a function of both the age of the child participating (x-axis) and the age of the child who originally produced the drawing (each line represents a different age). Error bars depict 95% bootstrapped confidence intervals.

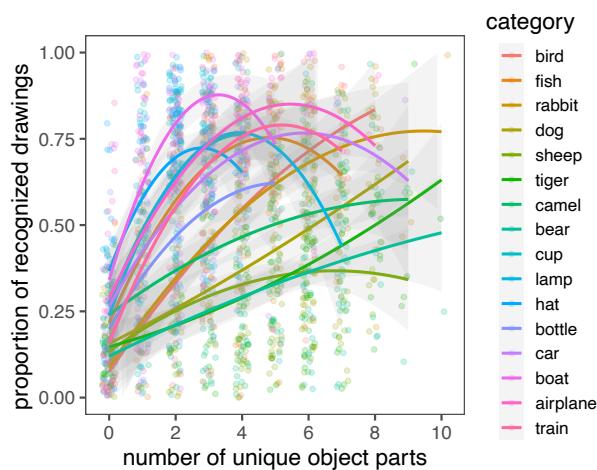
	Estimate	SE	z value	Pr(> z )
(Intercept)	0.05	0.23	0.20	0.84
Parts	129.39	10.04	12.89	<0.001
Parts**2	-34.98	3.24	-10.79	<0.001
Age	0.34	0.02	17.04	<0.001
Parts x Age	12.70	2.71	4.69	<0.001
Parts**2 x Age	-8.95	2.46	-3.64	<0.001

Table B3

*All model coefficients from a generalized, linear mixed effect model predicting how well children could recognize drawings of visual concepts as a function of their own age (Age; recognizer age) and the number of unique parts included in each drawing.*



*Figure B3.* Children's drawing recognition behavior for each of the 16 categories included in the recognition games; categories are grouped by the respective 4AFC game they were embedded in. Dots are scaled by the amount of data available from each age for each category (younger children were more frequent participants). Photo icons for each category are shown in the bottom right of each panel.



*Figure B4.* Drawing recognition for each category as a function of the number of unique parts included in each drawing; each individual dot is a unique drawing.