

**Hypothesis generation through active inductive inference in children and adults**

Neil R. Bramley\*

Department of Psychology, University of Edinburgh, Scotland

Gwyneth Heuser

Psychology Department, University of California, Berkeley, USA

Fei Xu

Psychology Department, University of California, Berkeley, USA

**Author Note**

Corresponding author\*: neil.bramley@ed.ac.uk.

Developmental data was collected under IRB protocol (Ref No: 2019-10-12687).

Adult data was collected under ethical approval granted by the Edinburgh University Psychology Research Ethics Committee (Ref No: 3231819/1). Supplementary material including all data and code is available at

[https://github.com/bramleyccslab/computational\\_constructivism](https://github.com/bramleyccslab/computational_constructivism). This study was not preregistered. Thanks to Jan-Philipp Fränken for help with coding free text responses. This research was supported by an EPSRC New Investigator Grant (EP/T033967/1) to N.R. Bramley and an NSF Award SMA-1640816 to F. Xu.

**Abstract**

A defining aspect of being human is an ability to reason about the world by generating and adapting ideas and hypotheses. Here we explore how this ability develops by comparing children’s and adults’ active search and hypothesis generation patterns in a task that mimics the open ended process of scientific induction. In our experiment, 54 children (aged  $8.97 \pm 1.11$ ) and 50 adults performed inductive inferences about a series of symbolic concepts through active testing. Children generated substantially more complex guesses about the hidden rule and were more elaborate in their testing behavior. We take a ‘computational constructivism’ perspective to explaining these patterns, positing that these inferences are driven by a combination of thinking (recombining and modifying existing concepts) and exploring (actively investigating and discovering patterns in the physical world). We show how our approach and rich new dataset help explain developmental differences in their hypothesis generation, active learning and inductive generalization.

## Hypothesis generation through active inductive inference in children and adults

1 A central question in the study of both human development and reasoning is how  
2 learners come up with new ideas and hypotheses to explain the world around them.  
3 Children excel at forming new categories, concepts, and causal theories (Carey, 2009) and  
4 by maturity, this coalesces into a unrivaled capacity for intelligent thought. Recent work  
5 in machine learning has begun to characterize computational principles driving both  
6 human development and mature human intelligence, pointing at a capacity to flexibly  
7 adapt, combine and re-purpose representations in order to generate new theories (Bramley,  
8 Dayan, Griffiths, & Lagnado, 2017; Bramley, Rothe, Tenenbaum, Xu, & Gureckis, 2018;  
9 Ellis et al., 2020); plans (Lai & Gershman, 2021; Lake, Ullman, Tenenbaum, & Gershman,  
10 2017; Ruis, Andreas, Baroni, Bouchacourt, & Lake, 2020); solve problems (Rule, Schulz,  
11 Piantadosi, & Tenenbaum, 2018) and create tools and technologies (Allen, Smith, &  
12 Tenenbaum, 2020). A related line of work has focused on understanding developmental  
13 change in the cognitive processes and behaviors that drive learning and underpin  
14 intelligence. The “child as scientist” (Gopnik, 1996) — or more recently, “child as hacker”  
15 (Rule, Tenenbaum, & Piantadosi, 2020) — perspective casts children’s cognition as driven  
16 by broadly the same processes as adults’ but at an earlier stage in a journey of  
17 construction and discovery. The idea is that childlike and adultlike behavioral differences  
18 should be reflected by parametrizable differences in search and thinking, reflecting rational  
19 principles of construction. Children’s hypothesis generation and search has also been  
20 framed as “higher temperature” than adults’ — producing more diversity of ideas at the  
21 cost of being noisier (Lucas, Bridgers, Griffiths, & Gopnik, 2014) — more narrowly focused  
22 on a few hypotheses at a time (Ruggeri & Lombrozo, 2014), and driven more by directed  
23 exploration and less by generalization (Wu, Schulz, Speekenbrink, Nelson, & Meder, 2017).

24 Constructivism is an influential perspective in developmental psychology (Carey,  
25 2009; Xu, 2019) and philosophy of science (Fedyk & Xu, 2018; Quine, 1969) that posits  
26 learners actively construct new ideas through a mixture of thinking—recombining and  
27 modifying ideas—and play—exploring and discovering patterns in the world (Bruner, Jolly,  
28 & Sylva, 1976; Piaget & Valsiner, 1930; Xu, 2019). While influential, constructivism has  
29 so-far lacked a formal model, limiting its testable predictions and ability to contribute to  
30 AI development. Indeed, constructivist ideas are almost absent from major branches of  
31 current computational cognitive science. For instance, Bayesian models have played a  
32 dominant role in recent study of cognition, providing a principled way of modeling  
33 probabilistic inference (Howson & Urbach, 2006). However, since Bayesian accounts  
34 describe learning within a predefined hypothesis space, they are silent about how a learner  
35 might explore or extend that space. Neural networks have also received much recent

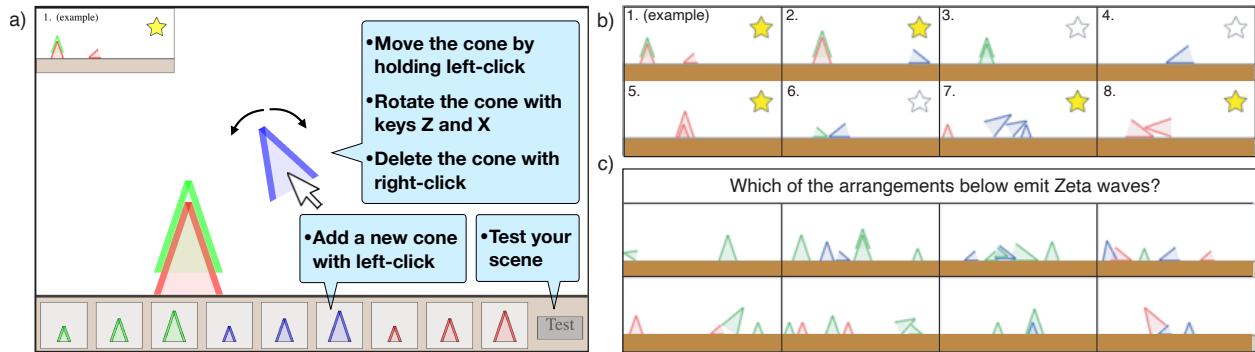
36 attention and achieve human-level performance in complex pattern recognition tasks  
37 (LeCun, Bengio, & Hinton, 2015), yet they require large amounts of training data  
38 compared to people and fail with unfamiliar inputs (Lake et al., 2017). From the  
39 constructivist view, this is to be expected. Neural network architectures are normally fixed  
40 ahead of training, and the weight-based representations they form are both opaque and  
41 distributed, limiting their suitability for later recombination. Information theory has also  
42 featured frequently in cognitive science as a metric of idealized information acquisition  
43 (Bramley, Dayan, et al., 2017; Gureckis & Markant, 2012). However, information-theoretic  
44 analyses also presuppose the Bayesian notion that learners have the relevant possibilities in  
45 mind and act to discriminate between them, rather than to support the task of constructing  
46 or discovering better ones. The central goal of this paper is develop a computational model  
47 of constructivist inference and use it as a lens to examine children's and adults' learning.

48 Traditionally, psychology research takes either a qualitative approach — studying  
49 cognition in naturalistic settings, eschewing formal models and statistics (Clarke & Braun,  
50 2014; Piaget & Valsiner, 1930) — or a quantitative approach — distilling aspects of  
51 learning into constrained tasks to allow statistical tests of theories (Pearson, 1930).  
52 However, the lure of quantitative measurement has led to narrow focus on tasks that miss  
53 essential aspects of real-world learning. For example, learning studies are usually explicit  
54 about the possible hypotheses, and limit actions and responses, making the task one of  
55 discrimination or choice rather than generation. Bayesian models provide sensible  
56 benchmarks for these settings (Anderson, 1990; Marr, 1982), but it is unclear what  
57 behavioral alignment with them reveals (M. Jones & Love, 2011), since in real  
58 environments, good hypotheses are hard won and behavior is limited only by imagination.  
59 While constrained tasks bring advantages of convenience, the worry is they “short circuit”  
60 cognition, blocking the active, deeply generative aspects of naturalistic learning (Gureckis  
61 & Markant, 2012).

62 Quantitative psychology traditionally minimizes interpretation of introspective  
63 self-report data (e.g., Dennett, 1991; Johansson, Hall, & Sikström, 2008), but see also  
64 (Newell & Shanks, 2014; Szollosi, Liang, Konstantinidis, Donkin, & Newell, 2019). Indeed,  
65 a dominant perspective views concepts as similarity-based clusters of features (Medin &  
66 Schaffer, 1978; Posner & Keele, 1968; Shepard & Chang, 1963) that drive categorization of  
67 new examples (Kruschke, 1992; Love, Medin, & Gureckis, 2004) but lack symbolic  
68 structure. However, people often not only generalize successfully but can also describe their  
69 concepts to others, combine them in imagination and draw analogies (Holyoak & Thagard,  
70 1989). It is not clear how subsymbolic approaches can account for these capabilities. We  
71 thus build on another line of work that considers concepts as symbolic and compositional

72 in character (Bramley et al., 2018; Goodman, Tenenbaum, Feldman, & Griffiths, 2008;  
 73 Piantadosi, Tenenbaum, & Goodman, 2016). We stake a new methodological path,  
 74 studying human hypothesis generation in an interactive context, where participants  
 75 generate their own evidence as they learn and provide both forced-choice generalizations  
 76 (quantitative data) and free guesses (qualitative data). We use a constructivist framework  
 77 to bridge the divide between the two forms of data, by providing a formal model of  
 78 hypothesis generation that can account for both quantitative and qualitative responses.  
 79 Thus, our approach provides insight into the source of human flexibility in learning and  
 80 thinking obscured by research focused on constrained tasks.

81 Our task is inspired by a tabletop game of scientific induction called “Zendo”  
 82 (Heath, 2004). In it, learners both observe and create *scenes*, which are arrangements of  
 83 2D triangular objects called *cones* (Figure 1) and test them to see if they produce a causal  
 84 effect. The goal is to predict which a set of new scenes will produce the effect and identify  
 85 the hidden rule that determines the general set of circumstances produce the effect (try it  
 86 [here](#)). Scenes could contain a varied number of cones. Each has two immutable properties:  
 87 size  $\in \{\text{small, medium, large}\}$  and color  $\in \{\text{red, green, blue}\}$  and continuous scene-specific  
 88  $x \in (0,8)$ ,  $y \in (0,6)$  positions and orientations  $\in (0,2\pi)$ . In addition to cones’ individual  
 89 properties, scenes also admit many relational properties arising from the relative features  
 90 and arrangement of different cones. For instance, subsets of cones might share a feature  
 91 value (i.e., be the same color, or have the same orientation) or be ordered on another (i.e.,  
 92 be larger than, or above) and pairs of cones might have relational properties like pointing  
 93 at one another or touching. This results in a rich implicit space of potential concepts.



**Figure 1**

The experimental task: a) Active learning phase. b) An example sequence of 8 tests, the first is provided to all participants, and subsequent tests are constructed by the learner using the interface in (a). Yellow stars indicate those that follow the hidden rule. c) Generalization phase: Participants select which of a set of new scenes are rule following by clicking on them.

94 In order to model the task, we adopt an expressive concept grammar inspired by  
 95 “constructivist” ideas in developmental psychology and formalized using “program  
 96 induction” ideas from machine learning. Concretely, we assume the latent space of possible  
 97 concepts in our task are those expressible in first order logic combined with lambda  
 98 abstraction and full knowledge of the potentially relevant features of the scene (see  
 99 Appendix Table A-1 for the grammatical primitives we assume). Table 1 shows the five  
 100 ground truth rules we used in our experiment expressed in natural language and in lambda  
 101 calculus along with the initial rule-following example scene we provided to participants.

## 102 Context-free hypothesis generation

103 In accounting for children’s and adults’ inferences, we entertain two related  
 104 constructivist algorithms. The first takes a fully “top down” approach to inference,  
 105 utilizing a probabilistic context-free grammar (PCFG) to define a latent prior over  
 106 concepts expressible in first order logic. A PCFG is a collection of “productions” that  
 107 stochastically build expressions in an underlying grammar (Ginsburg, 1966). A PCFG can  
 108 be used to generate a prior sample of hypotheses that can then be weighted by their  
 109 likelihoods of producing observations—here, their ability to reproduce the labels for the  
 110 scenes that the participant has tested. The model’s best guess about the hidden rule is  
 111 then the *maximum a posteriori* hypothesis in the sample. The hypotheses make predictions  
 112 about new scenes which can be weighted by their posterior probability and marginalized  
 113 over to make generalizations. Because parts of the production process and underlying  
 114 grammar involve branching—e.g., “and” and “or”—hypotheses can become arbitrarily long  
 115 and complex, involving multiple Boolean functions and complex relationships between an  
 116 unlimited number of bound variables. In this way, an infinite latent space is covered in the  
 117 limit of infinite PCFG sampling (see Figure 2a).

118 The probabilities for each production in a PCFG can be fit to maximize  
 119 correspondence with human judgments. Different PCFGs, containing different primitives  
 120 and expansions, can be compared against human behavior. In this way, recent work has  
 121 attempted to infer the “logical primitives of thought” (Goodman et al., 2008; Piantadosi et  
 122 al., 2016). In the current work we consider a single expressive PCFG architecture but  
 123 contrast its behavior uniform production weights with its behavior with fitted “childlike”  
 124 and “adultlike” weights, that allow it to reproduce the summary statistics of children’s and  
 125 adults’ guessed rules. Crucially, under all three weightings, our PCFG embodies the  
 126 principle of parsimony: Simpler concepts—composed of fewer grammatical parts (Feldman,  
 127 2000)—have a higher prior probability of being produced and so are favored over more  
 128 complex ones equally able to explain the data.

What PCFG approaches have in common is a generative mechanism for sampling from an infinite latent prior, here over possible logical concepts. However, sampled “guesses” must then be tested against data. Unfortunately, most samples are likely to be inconsistent with whatever data a learner has already encountered. For this reason, the procedure is inherently inefficient, and requires a very large numbers of samples in order to reliably identify non-trivial rules. Thus, we also consider an alternative that provide a more computationally plausible inference mechanism.

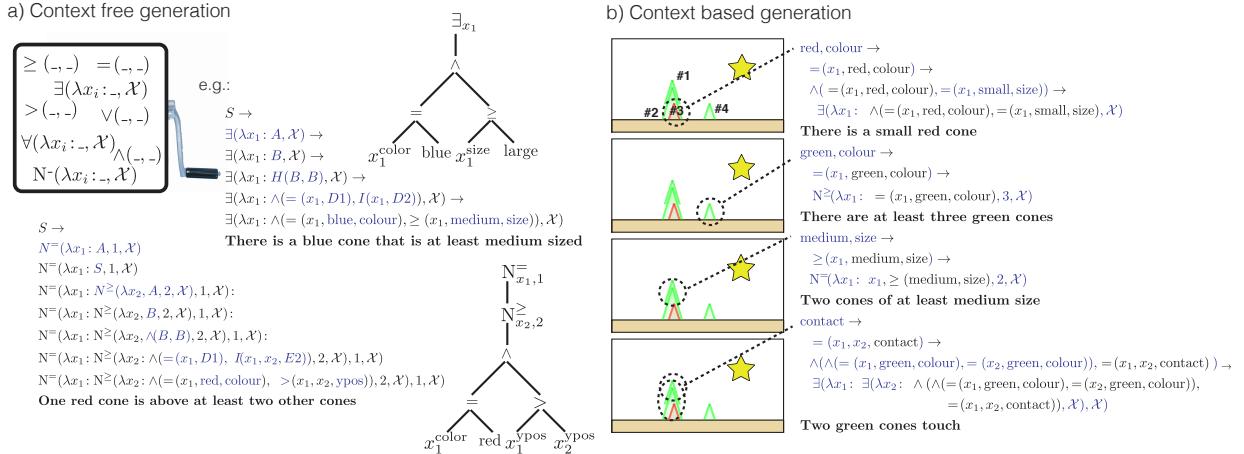
### **136 Context-based hypothesis generation**

Instance Driven Generation (IDG) (Bramley et al., 2018) is a recent proposal related to the PCFG but with one key difference. Rather than generating hypotheses *a priori*, it generates ideas *inspired* by encountered examples (cf. Michalski, 1969), thus blending top-down generation with bottom-up reactivity to evidence. An IDG learner starts by observing the features of objects in a scene and uses these to back out a true logical statement about the scene in a stochastic but truth-preserving way. If the scene is rule following, this statement constitutes a positive hypothesis about the hidden rule. Otherwise, it constitutes a negative hypothesis, i.e. about what must *not* be present. Thus, IDG does not generate uniformly from all possible concepts, but directly from a restricted space consistent with a focal observation. Figure 2b illustrates this approach. While the PCFG starts at the outside and works inward, the IDG starts from the central content and works outward out to a quantified statement, ensuring at each step that it is true of the scene. As with the PCFG, we consider a uniform variant as well as variants that include productions reverse engineered to match the summary statistics of guesses generated by children and by adults.

### **152 Active learning**

Children have long been seen as primarily active learners, using “play” to explore their environment and test their hypotheses (Bruner et al., 1976; Cook, Goodman, & Schulz, 2011; Piaget & Valsiner, 1930). Information theory is used to benchmark active learning (Nelson, 2005; Shannon, 1951) but assumes learners know the relevant possibilities and acts to discriminate rather than to constructing or discovering better ideas, potentially explaining why behavioral alignment with information maximization is mixed and task dependent (Coenen, Nelson, & Gureckis, 2018). The developmental literature has emphasized the utility of “control of variables” heuristic (Chen & Klahr, 1999; A. Jones, Bramley, Gureckis, & Ruggeri, in revision; Klahr, Fay, & Dunbar, 1993; Klahr, Zimmerman, & Jirout, 2011) — manipulating one design variable at a time, so that

163 changes in the outcome can be unambiguously attributed to the change in the input. Past  
 164 research has only focused on restricted settings with a few simple variables and our task is  
 165 much more complex. Thus, in exploring the active learning in our task, we will look for the  
 166 empirical signature of control of variables style incremental and systematic testing.

**Figure 2**

a) Example generation of hypotheses using the PCFG. b) Examples of IDG hypothesis generation based on an observation of a scene that follows the rule. New additions on each line are marked in blue. Full details in Appendix.

167 In sum, the core contribution of this work is a close investigation of developmental  
 168 differences in active open ended hypothesis generation and development of constructivist  
 169 modeling approach that bridges qualitative–quantitative and symbolic–subsymbolic  
 170 divides. To foreshadow, we find evidence of compositional concept formation in both adults  
 171 and children; support for a bottom-up instance driven account. We find children create  
 172 more complex learning data although do so less systematically than adults. They then go  
 173 on to make more complex guesses while achieving a commensurate fit to evidence. Our  
 174 constructivist framework suggests this behavior is a natural result of “flatter” idea and  
 175 action generation mechanisms.

176

## Experiment

177 **Methods**178 **Participants**

179 We recruited 54 children in the lab (23 female, aged  $8.97 \pm 1.11$ ) and 50 adults  
 180 online (22 female, aged  $38.6 \pm 10.2$ ). Forty children completed all five trials and the  
 181 remaining 14 completed  $2.71 \pm 1.07$  trials before indicating that they had had enough. For

**Table 1**  
*Rules Tested in Experiment*

Rule	Initial Example
1. There's a red $\exists(x_1: = (x_1, \text{red}, \text{color}), \mathcal{X})$	
2. They're all the same size $\forall(x_1: \forall(x_2: = (x_1, x_2, \text{size}), \mathcal{X}), \mathcal{X})$	
3. Nothing is upright $\forall(x_1: \neg(= (x_1, \text{upright}, \text{orientation})), \mathcal{X})$	
4. There is exactly 1 blue $N=(\lambda x_1: = (x_1, \text{blue}, \text{color}), 1, \mathcal{X})$	
5. There's something blue and small $\exists(x_1: \wedge(= (x_1, \text{blue}, \text{color}), = (x_1, 1, \text{size}), \mathcal{X})$	

<sup>182</sup> these children we simply include the trials that they completed. We collected participants  
<sup>183</sup> until we reached our intended sample size of 50 per agegroup after exclusions. Ten  
<sup>184</sup> additional adult participants completed the task but were excluded before analysis for  
<sup>185</sup> providing nonsensical or copy-pasted text responses. Adult participants were paid \$1.50  
<sup>186</sup> and as well as a performance related bonus of up to \$4 ( $\$1.96 \pm 0.75$ ). For children sessions  
<sup>187</sup> lasted between 30 minutes and an hour. For adults, the task took  $27.49 \pm 12.09$  minutes of  
<sup>188</sup> which  $9.8 \pm 7.9$  was spent on instructions.

### <sup>189</sup> Materials and Procedure

#### <sup>190</sup> Child sample.

<sup>191</sup> **Instructions.** Participants sat in front of a laptop with a mouse attached, with

<sup>192</sup> the experimenter sitting next to them. The laptop displayed this webpage:

<sup>193</sup> [http://www.bramleylab.ppls.ed.ac.uk/experiments/zendo\\_kids/task.html](http://www.bramleylab.ppls.ed.ac.uk/experiments/zendo_kids/task.html).

<sup>194</sup> The experimenter read out the instructions displayed on the webpage for the  
<sup>195</sup> participant. These explained how the game worked and showed the participant five  
<sup>196</sup> examples of possible rules the blocks could have (relating to color, size, proximity, angle, or  
<sup>197</sup> relation). The instructions also included videos showing the participant how to manipulate  
<sup>198</sup> the blocks using the mouse and keyboard. After the instructions, the participant was given  
<sup>199</sup> a comprehension check of five true or false questions. If they did not get them all right on  
<sup>200</sup> their first try, the experimenter read through the instructions again and asked them again.

201 All participants passed the comprehension check the second time. The participant was  
202 then introduced to an initial example of a block type (“Here are some blocks called  
203 [name]s. We’re going to click test to see if stars will come out of the [name]s.”). The initial  
204 example of each block type (i.e., each rule) was constant across participants. There were  
205 five block types in total, one for each rule, and the order of these was randomized (see  
206 Table 1). Every initial example of a block type was a positive example, so a star animation  
207 played when the “Test” button was clicked. The participant was encouraged to use either  
208 the trackpad or the mouse to click the “Test” button, whichever was comfortable for them.

209 **Learning Phase.** After the initial positive example, the participant was shown a  
210 blank scene with blocks available to add to it, and was asked to test the blocks seven more  
211 times (Figure 1a). The scene creation interface was subject to simulated gravity, meaning  
212 there were physical constraints on how the objects can be arranged. The experimenter told  
213 them they could now play with the blocks like they saw in the instructional video. The  
214 experimenter also reminded the participant of how to add, remove, move, and rotate blocks  
215 on the screen using the mouse and keyboard. Participants were encouraged to ask for help  
216 with moving the blocks if needed. If they seemed to be having trouble, the experimenter  
217 would ask if they needed help with setting up the blocks. The participants were told that  
218 when they were done moving the blocks around, they should press the “Test” button to see  
219 if stars came out of them. For positive tests, the experimenter would neutrally say: “Stars  
220 did come out of the [name]s that time” and for negative tests: “Stars did not come out of  
221 the [name]s that time.”

222 **Question Phase.** After testing the blocks a total of eight times ((Figure 1b)),  
223 participants were shown a selection of eight more pre-determined scenes containing blocks  
224 (Figure 1c). The experimenter asked them to click on which pictures they thought the  
225 stars would come out of, reminding them that they could pick as many as they wanted, but  
226 they had to pick at least one. Unknown to participants, half of these scenes were always  
227 rule following but their positions on screen were independently counterbalanced.

228 **Free Responses.** Participants were then presented with a blank text box and  
229 asked, “What do you think the rule is for how the [name]s work?” The experimenter typed  
230 into the text box the participant’s verbal answer verbatim, or as close as possible.

231 The Testing, Question, and Free Response phases were repeated identically for each  
232 of the five block types. After the five trials were completed, the participant was shown the  
233 results including each true rule and how well they did on each problem and was thanked for  
234 playing the game. As compensation, participants were allowed to pick a small toy out of a  
235 prize box, and parents were given a paper “diploma” to commemorate their child’s visit.

236 **Adult sample.** We recruited our adult sample from Amazon Mechanical Turk

237 and adults completed the task on their own computers. They completed the same  
238 instructions as the children with an additional section about bonuses and had to  
239 successfully answer comprehension questions, including an additional two about the  
240 bonuses, before starting the main task. Specifically, adults were bonused 5 cents for each  
241 correct generalization (up to a possible 40 cents for each of the five trials) and an  
242 additional 40 cents for an correct guess as to the hidden rule, again for each of the five  
243 trials. Aside from having no experimenter in the room, and filling out the text fields  
244 themselves, the procedure was identical to the children’s task. Full materials including  
245 experiment demos, data and code are available at the [Online Repository](#).

246 **Results**

247 We first look at the qualitative characteristics of children’s and adults’ explicit rule  
248 guesses then assess relative accuracy of participants’ rules and generalizations about new  
249 scenes. We compare children’s accuracy to adults’ and both to our constructivist learning  
250 algorithms: Fully top down context-free generation from an expressive latent prior —  
251 Probabilistic Context Free Generation (PCFG) — and a partially bottom-up generation —  
252 Instance Driven Generation (IDG) (Bramley et al., 2018). We then turn to analysis of the  
253 scenes produced by adults and children and finally evaluate a set of formal models’ ability  
254 to produce both participants generalizations and their encoded free responses.

255 **Rule complexity and constituents**

256 We had human coders translate participants’ free text guesses about the hidden rule  
257 wherever possible into equivalent logical lambda expression using the grammatical elements  
258 available to our learning models. We were able to do this for 86% ( $n=205$ ) of children’s  
259 trials and 88% ( $n=219$ ) of adults’ trials. For example, if the participant wrote “*There must*  
260 *be one big red block*” this was converted into

$$\text{261 } N = (\lambda x_1 : \wedge (= (x_1, \text{large}, \text{size}), = (x_1, \text{red}, \text{color})), 1, \mathcal{X}).^1$$

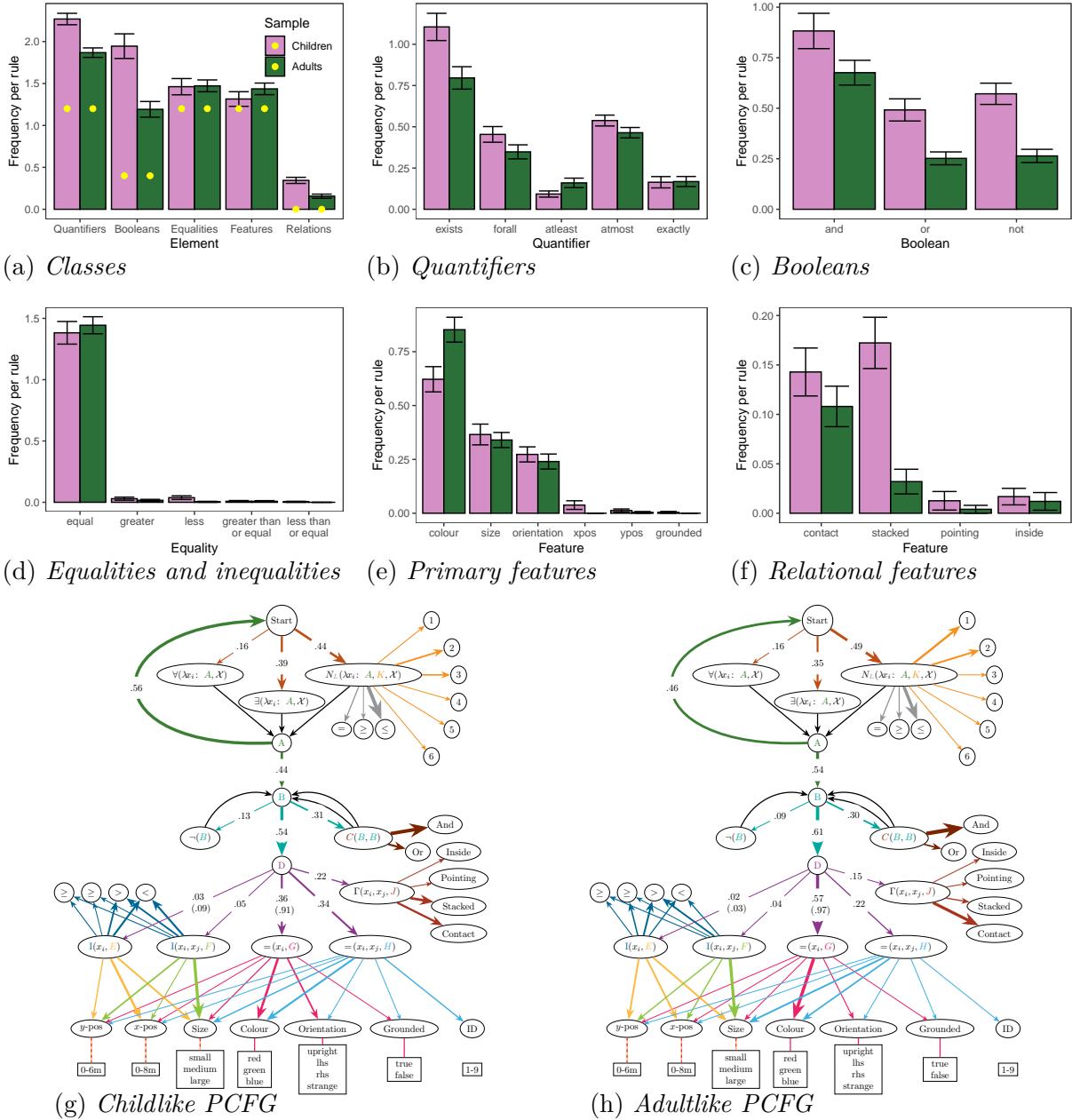
262 To explore structural differences in children’s versus adults’ hypotheses, we first  
263 break down these encoded rule guesses into their grammatical parts. This reveals that  
264 children’s encoded rules were substantially more complex than those generated by adults  
265 and that both were more complex than the ground truth rules. Children’s and adults’ rules

---

1 This logical version can be automatically evaluated on the scenes and can be read literally as asserting “*There exists exactly one  $x_1$  in the set of objects  $\mathcal{X}$  such that  $x_1$  has the size ‘large’ and the color ‘red’*”. We provide details about the coding in the Appendix and full coding resources and full coding data in the Online Repository.

266 also differed in terms of the prevalence of particular elements and features (see Figure 3).  
 267 As an example, one child’s rule was “*You must have two reds and one blue*” which was  
 268 translated to  $N^=(\lambda x_1: N^=(\lambda x_2: (\wedge(=(x_1, \text{red}, \text{color}), =(x_2, \text{blue}, \text{color})), 1, \mathcal{X}), 2, \mathcal{X})$ ,  
 269 requiring two quantifiers ( $N^=$ ), one boolean ( $\wedge$ ), 2 equalities ( $=()$ ), and two references to  
 270 the feature color. The typical child-generated-rule used 2.25 quantifiers (3b), 2.06 booleans  
 271 (3c), 1.55 equalities and inequalities (3d), referred to 1.39 different primary features (color,  
 272 size, orientation, x- or y-position, groundedness, 3e) and 0.37 relational features (contact,  
 273 stackedness, pointing, or insideness, 3f). In contrast, the average adult generated rule  
 274 required just 1.84 quantifiers, 1.20 booleans, 1.47 equalities and inequalities, and referred  
 275 to 1.44 primary features but only 0.16 relational features. When children posited that an  
 276 “at least”, “at most” or “exactly” a certain number of objects must have certain features,  
 277 the number they chose were substantially higher than that for adults (2.36 compared to  
 278 1.58). In terms of features, adults strongly tended to posit rules relating to color (58%  
 279 compared to 39% of children’s rules), while children were more likely to refer to positional  
 280 properties (26% compared to 18% of adults’ rules) and relations (31% compared to 14% of  
 281 adults’ rules) between the objects.

282 **Reverse engineering Childlike and Adultlike prior productions.** Having  
 283 encoded all the rule guesses from adults and children, we can work back from the  
 284 distribution of rules to create a set of productions that produces a similar sample. To do  
 285 this we work back from the observed counts for each rule element. To roughly  
 286 accommodate the fact that rules are conditional on different data, we regularized these  
 287 counts by including a prior pseudo-count of 5 on all productions. For example, children’s  
 288 rules involved  $\exists$  263 times,  $\forall$  108 times and  $N$  297 times, so we assumed prior production  
 289 weights of  $\{263 + 5, 108 + 5, 297 + 5\}/(263 + 108 + 297 + 15) = \{.39, .17, .44\}$ . The full set  
 290 of fitted weights for both adults and children are visualized and detailed in Figure 3g&h.  
 291 Strictly these are samples from a range of different participants’ posteriors  $P(r|d_{p,t})$  not  
 292 from their prior  $P(r)$ , since judgments were always conditional on some evidence. However,  
 293 since evidence differs greatly across the rules we considered and scenes participants created,  
 294 and since the structural elements of the grammar (Booleans, Quantifiers etc) are not  
 295 tightly tied to scene-specifics, we feel this still provides an informative and useful  
 296 elucidation of differences in a common set of productions that can produce children and  
 297 adults’ hypotheses. This analysis illustrates that children’s production process is “flatter”  
 298 than adults’ under a constructive account, with a greater average entropy over the various  
 299 production steps of this process  $1.28 \pm 0.50\text{bits}$   $1.03 \pm 0.59\text{bits}$ ,  $t(13) = 3.2, p = 0.007$ .

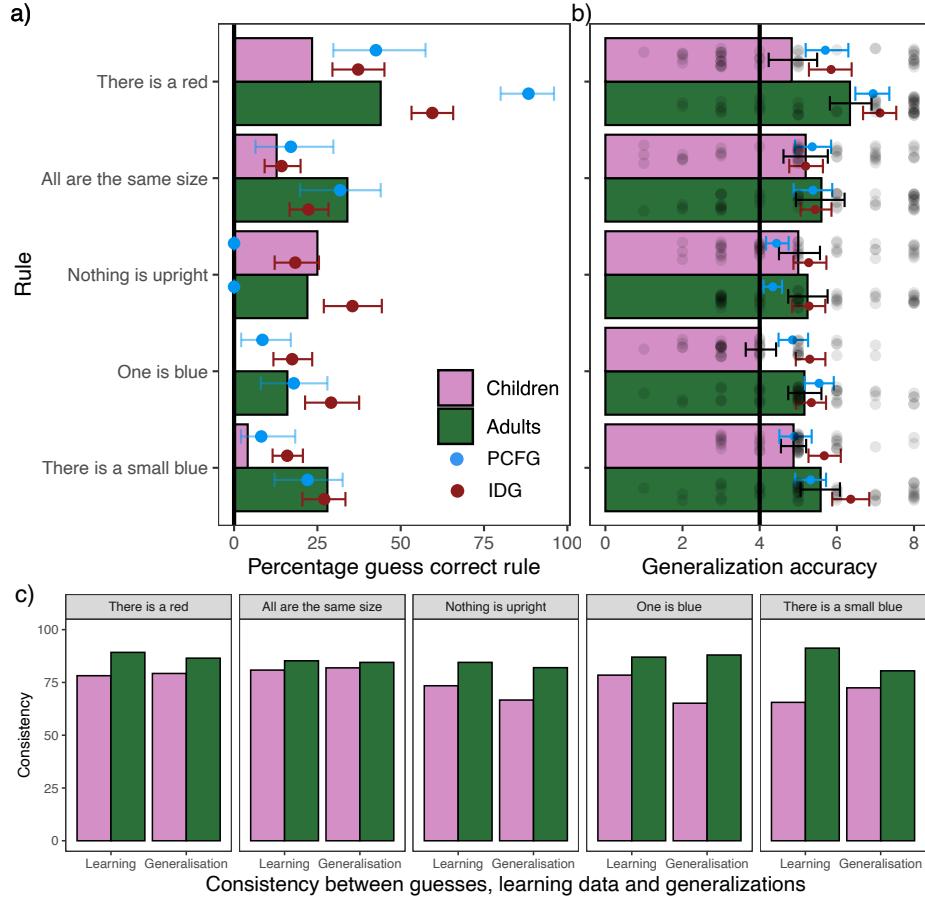
**Figure 3**

(a-f) Relative frequency of rule elements in Children's and Adults' rule guesses. Yellow points show ground truth frequencies. (g&h) Visualization of the childlike and adultlike PCFGs, reverse engineered to produce rules with empirical frequencies matched to children's and adults'. A rule is produced by following arrows from "Start" according to their probabilities (line weights and annotation), replacing the capital letters with the syntax fragment at the arrow's target and repeating until termination.

### 300 Accuracy

301 Having observed systematic differences in the content of children's and adults'  
 302 hypotheses, we now ask if these differences manifest in children's and adults' inferential

303 success; their ability to identify the ground truth and make accurate generalizations.



**Figure 4**

a) Percentage of participants guessing exactly the correct rule. Bars show mean  $\pm$  bootstrapped 95% confidence intervals for children (pink) and adults (green). b) Generalization performance. Blue and red points show mean performance of PCFG and IDG simulations with fitted production weights  $\pm$  bootstrapped 95% confidence intervals. Black vertical lines denote chance performance. c) Consistency between subjects' rule guess and their (self generated) learning data, and generalization judgments.

304 **Rule guesses.** Both children and adults were sometimes able to guess exactly the  
 305 correct rules, doing so a respective 11% and 28% of trials. Adults produced the correct rule  
 306 more frequently than children  $t(102) = 4.0, p < .001$  and were more likely than children to  
 307 guess correctly (at a corrected significance level of 0.01) for the "All are the same size",  
 308 "One is blue" and "There is a small blue" rules (see Figure 4a). Note that chance level  
 309 baseline for this kind of guess is essentially 0%. There are an unlimited number of wrong  
 310 guesses and a small set of semantically correct guesses. It is also the nature of this  
 311 inductive problem that there is an infinite number of perfectly consistent rules for any  
 312 evidence, although as more evidence arrives the ground truth is increasingly likely to be

313 among “simplest” rules in this set. Thus, it is instructive to ask whether participants rules,  
 314 where not exactly correct, are still consistent with the evidence they have seen.

315 Children’s explicit rule guesses were consistent with the labels of all 8 training  
 316 scenes 30% of the time while Adult’s guesses were fully consistent 54% of the time. A  
 317 completely random rule would only be consistent with all 8 scenes around  
 318  $0.5^8 \times 100 = 0.4\%$  of the time. There was a moderate difference in average proportion of  
 319 learning data explained by children’s compared to adult’s rules  $71\% \pm 27\%$  vs  $87\% \pm 17\%$   
 320  $t(98) = 5.6, p < .001$ . Similarly there was a difference the proportion of the participants’  
 321 generalizations that were consistent with their rule guess  $72\% \pm 21\%$  vs  $84\% \pm 16\%$ ,  
 322  $t(98) = 4.1, p < .001$  (see Figure 4c for a by-rule breakdown).

We now compare this to simulated *context free* (PCFG) and *context based* (IDG) learning algorithms provided with the active learning data generated by the human participants. We produced 50,000 childlike rules  $\hat{R}_c$  and 50,000 adultlike rules  $\hat{R}_a$  that have properties matched to those in Figure 3a–f as well as an additional sample of rules based on uniform production weights  $\hat{R}_u$ . These act as an approximation to the infinite latent prior over rules  $P(r)$ , before seeing any data. In order to approximate a posterior over rules given self-generated learning scenes  $\mathbf{d}$ , we then weight these samples by their likelihood of producing the scene labels observed during the learning phase

$$P(r|\mathbf{d}) \propto P(\mathbf{d}|r)P(r) \quad (1)$$

$$\approx P(\mathbf{d}|r) \sum_{\hat{r} \in \hat{R}} \mathbb{I}(r = \hat{r}) \quad (2)$$

323 and count how often they appear in the prior sample, with indicator function  $\mathbb{I}(\cdot)$  denoting  
 324 exact or semantic equivalence. To test for semantic equivalence, we computed predictions  
 325 for the first 500 participant-generated scenes for each rule and clustered together those that  
 326 made identical predictions. We round positional features to one decimal place in evaluating  
 327 rules to allow for perceptual uncertainty.

328 We assume the following likelihood function

$$P(\mathbf{d}|r) = \exp(-b \times N_{\text{outliers}}) \quad (3)$$

329 embodying the idea that: the more training points a rule cannot explain, the less  
 330 likely it is to be true. For a large  $b$  the likelihood function approaches the true  
 331 deterministic behavior of the rules. However, to allow for some noise while maintaining  
 332 computational tractability in our analyses we simply assume a  $b = 2$ , corresponding to a  
 333 likelihood function that decays rapidly from 1 for rules that predict all 8 scenes’ labels to

<sup>334</sup> .13 for a single misprediction, and .02 for 2 mispredictions and so on.

<sup>335</sup> To generate IDG predictions, as with the PCFG, we produced a childlike, an  
<sup>336</sup> adultlike and a uniform sample of instance driven hypotheses. This involved merging the  
<sup>337</sup> production probabilities from the PCFG into the Instance Driven Generation procedure  
<sup>338</sup> detailed in the Appendix. However, since each generation depended on the particular  
<sup>339</sup> scenes for inspiration and this differed for every participant and trial, we generated smaller  
<sup>340</sup> samples of 10,000 childlike, adultlike, and uniform rules for each trial. We spread these  
<sup>341</sup> evenly across the 8 learning scenes. For scenes that did not follow the rule we followed the  
<sup>342</sup> same procedure as for scenes that did, but wrapped the rule in a negation. For example,  
<sup>343</sup> observing a non-rule-following scene in which there are objects in contact might inspire the  
<sup>344</sup> rule that no cones are touching.

<sup>345</sup> Taking the maximum *a posteriori* estimate (guessing in the event of ties) under  
<sup>346</sup> either model leads to guessing the correct hypothesis at similar levels to participants. For a  
<sup>347</sup> uniform-weighted PCFG sample, the MAP is correct on  $9\% \pm 28\%$  of children's trials and  
<sup>348</sup>  $12\% \pm 33\%$  of adults' trials. Note that since these simulations use the same prior sample,  
<sup>349</sup> the small differences we see are due to the different learning data generated by children and  
<sup>350</sup> adults. However, accuracy improves substantially and better reproduces the empirical  
<sup>351</sup> child-adult accuracy difference if we use samples based on reverse engineered weights that  
<sup>352</sup> reproduce the qualitative properties of children's and adults' rules (see Appendix and  
<sup>353</sup> Figures 3g&h). For these samples, we get correct PCFG guesses for  $15\% \pm 36\%$  of children's  
<sup>354</sup> trials and  $32\% \pm 46\%$  of adults' trials. Across rules, the PCFG does not match well with  
<sup>355</sup> children's or adults' accuracy, overperforming participants on the syntactically simpler rule  
<sup>356</sup> "there is a red" but failing to capture participants correct guesses of "Nothing is upright".

<sup>357</sup> Uniform-weight IDG simulations guess correctly on  $15\% \pm 18\%$  of trials for children's  
<sup>358</sup> learning data and  $25\% \pm 23\%$  for adults learning data. Using the reverse-engineered  
<sup>359</sup> weights, this increases to  $21\% \pm 23\%$  and  $35\% \pm 29\%$  and again provides a visually closer  
<sup>360</sup> fit to the by-rule guess rates (Figure 4a).

<sup>361</sup> **Generalizations.** We now analyze the quantitative response data constituted by  
<sup>362</sup> forced choice generalizations about which of 8 new scenes will produce stars (i.e. follow the  
<sup>363</sup> hidden rule). Across the five tasks, both children and adults guessed more accurately than  
<sup>364</sup> chance (50%): *children* mean $\pm SD$   $59\% \pm 11\%$ ,  $t(53) = 5.9, p < .001$ ; *adults*  
<sup>365</sup>  $70\% \pm 14\%$ ,  $t(49) = 10.3, p < .001$ . Adults' generalizations were significantly more accurate  
<sup>366</sup> than children's  $t(102) = 4.6, p < .001$  and children's accuracy improved significantly with  
<sup>367</sup> age  $F(1, 52) = 6.2, \eta^2 = .11, p = 0.015$ . Indeed, adults' generalization accuracy was above a  
<sup>368</sup> Bonferroni-corrected chance level of  $p \leq 0.01$  for all five rules and children were similarly  
<sup>369</sup> above chance except for rules 1. "There is a red" ( $t(46) = 2.5, p = .015$ ) and 4. "One is

370 blue” ( $t(46) = .1, p = .915$ ; see Figure 4b).

We compare this pattern against simulated constructivist PCFG and IDG learner benchmarks. To do this we use the requisite predictive distribution to model generalizations to the set of test scenes  $\mathbf{d}^*$

$$P(\mathbf{d}^*|\mathbf{d}) = \int_R P(\mathbf{d}^*|R)P(R|\mathbf{d}) dR \quad (4)$$

$$\approx \sum_{r \in \hat{R}} P(\mathbf{d}^*|r)P(r|\mathbf{d}) \quad (5)$$

371 Provided with the active learning data generated by the human participants, both  
 372 performed in the human range at generalization. Using uniform production weights and  
 373 taking the marginally most likely generalization labels over a posterior weighted sample of  
 374 PCFG-generated rules based on the participants active learning data yielded accuracies of  
 375  $60.3\% \pm 18.8\%$  for children’s and  $61.7\% \pm 19.5\%$  for adults’ data. The fitted-weight PCFG  
 376 models perform a little better and reproduces the empirical difference between children’s  
 377 and adults’ accuracy:  $63 \pm 20\%$  for children’s and  $69 \pm 21\%$  for adults’ PCFG weights. The  
 378 unfitted IDG, again, performed slightly better than the PCFG, generalizing at  
 379  $66.3\% \pm 20.1\%$  from children’s active learning data and  $69.5\% \pm 20.7\%$  from adults’. Again,  
 380 the fitted-weight IDG models’ performance increased slightly and better reproduced the  
 381 difference between children’s and adults’ accuracy ( $68 \pm 20\%$  and  $74 \pm 21\%$ ).

382 The better accuracy of the IDG compared to the PCFG replicates the findings of  
 383 (Bramley et al., 2018) and extends it to children as well as adults. Intuitively, this is  
 384 because the bottom-up mechanism ties the hypotheses generated to features of the learning  
 385 cases, effectively narrowing in on plausible hypotheses more efficiently. More broadly, these  
 386 simulation results underscore the inherent difficulty of inductive inference. Even in this  
 387 “small world” with known and fully observed features and allowing cognitively implausibly  
 388 large numbers of hypothesis samples, it is not possible to robustly outperform human  
 389 adults in this task. The PCFG and IDG were not statistically better or worse than  
 390 participants at any rule inference under after Bonferroni correction with the exception that  
 391 the IDG outperformed children on rule 4  $t(96) = 4.5, p < .0001$ .

### 392 *Interim discussion*

393 Children were only moderately less able to guess rules that fit the evidence than  
 394 adults and there were only moderate differences in the compatibility between children’s  
 395 and adults’ rules and their generalizations. However, children did overfit the evidence  
 396 more, essentially producing more complex and naïve characterizations of the rule-following  
 397 scenes than did adults. This can be seen in the larger number of quantifiers and relations

398 mentioned in children’s rules than in adults’.

399 A complicating factor is that children generated different data to adults. However,  
 400 our PCFG and IDG simulations suggest exposure to different data cannot explain most of  
 401 the accuracy differences between children and adults. Using identical production weights  
 402 and the scenes generated by adults and children led to only small differences in accuracy  
 403 for the PCFG and moderate for the IDG, while using a “flatter” set of productions fit to  
 404 match childlike rules, and a “sharper” set fit to adults’ rules, better reproduces the  
 405 accuracy patterns. We take this to suggest hypothesis generation differences are driving a  
 406 large portion of the differences in children’s and adult’s inductive inferences.

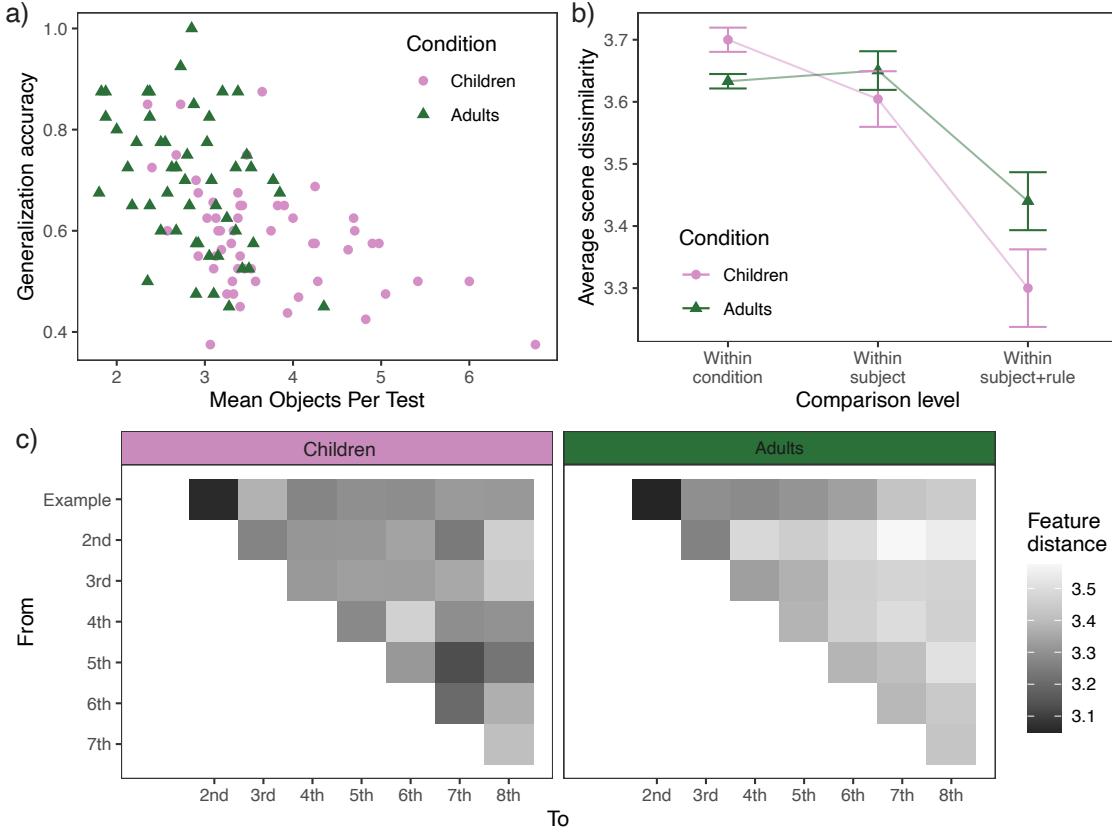
407 We now turn to analyze active learning (scene generation) behavior. We first  
 408 characterize the differences between the scenes generated by children and by adults and  
 409 then ask whether these can be attributed to differences in hypothesis generation.

410 ***Search behavior***

411 As well as generating more complex rules, children also tended to create  
 412 substantially more complex scenes than adults during the learning phase. The average  
 413 child-generated scene contained  $3.7 \pm 0.88$  objects compared to  $2.84 \pm 0.57$  objects for adults  
 414 ( $t(102) = 5.8, p < .001$ ). The complexity of test scenes was inversely related to performance  
 415 overall ( $F(1, 102) = 39.0, \beta = -0.08, \eta^2 = .28, p < .001$ ) and also within both the child  
 416 sample ( $F(1, 52) = , \beta = -0.056, \eta^2 = .20, p < .001$ ) and adult sample  
 417 ( $F(1, 49) = 9.1, \beta = -0.096, \eta^2 = .16, p < .001$ ) taken individually (see Figure 5a). Within  
 418 the child sample, age was inversely associated with scene complexity with an average of 0.35  
 419 fewer objects per scene for each additional year of age  $F(1, 52) = 12.6, \eta^2 = .19, p < .001$ .  
 420 Aside from this difference, we can also assess whether children’s or adults’ scenes bear the  
 421 hallmarks of a local “search” across possible scene dimensions.

422 **Scene sequences and similarity.** While we do not yet have a model of scene  
 423 creation process, we hypothesized that a *control of variables* strategy (Kuhn & Brannock,  
 424 1977) is a reasonable marker of systematic active learning. In the current setting, this  
 425 manifests as a tendency to generate new evidence by recreating a previous scene (i.e.  
 426 whose labels is already known) and making some change to it. This allows a learner to  
 427 isolate boundary conditions for the hidden rule, and so potentially fine tune a focal  
 428 hypothesis or rule among a small set of similar alternatives (Bramley, Dayan, et al., 2017).  
 429 Additionally, we speculated that reuse in general is likely to reduce cognitive load  
 430 (Gershman & Niv, 2010).

431 If this is the case, we should expect the scenes generated by participants to be more  
 432 similar to the initial example than to a random scene or scene drawn from a different

**Figure 5**

(a) Generalization accuracy by number of objects per test scene. (b) Average dissimilarity between self-generated scenes at different levels of aggregation. Error bars show standard errors for subject means. (c) Average similarity matrices between initial example and self generated scenes 2 to 8. See Appendix for detailed procedure and similarity matrices separated by component.

learning problem. To explore this, we constructed a distance metric that we used to measure the intuitive dissimilarity between any pair of scenes. The metric captures a form of edit distance, encoding how much and how many of the features (positions, colors, shapes) of the objects in one scene would have to be changed to reproduce the other scene. Essentially, this involved  $z$ -scoring and combining a “minimal-edit set” of feature differences and incorporating a proportional cost for additional or omitted objects. We provide a detailed procedure and example of how we computed these edit distances and break them down into their separate components in the Appendix. The mean distance between any randomly selected pair of participant-generated scenes was  $M \pm SD = 3.67 \pm 0.94$ . Taken as a whole, the scenes generated by children were more diverse than adults’ with average dissimilarity of  $3.70 \pm 0.14$  compared to  $3.63 \pm 0.08$ ,  $t(102) = 2.9, p = 0.0048$ .

However, this diversity seems to be *between* rather than *within* subject for children’s

445 choices. Within subject but across trials, the average inter-scene dissimilarity for children  
 446 was  $3.60 \pm .33$  similar to that for adults'  $3.65 \pm .22$ ,  $t(102) = .83, p = .4$ . Focusing more  
 447 narrowly, within the scenes produced by an individual subject while learning about a single  
 448 rule, we see a reversal of the aggregate pattern. That is, within a learning task, children's  
 449 scenes are marginally *less* diverse on average than adults' (children:  $3.30 \pm 0.459$ , adults:  
 450  $3.44 \pm 0.33$ ,  $t(102) = 1.77, p = 0.08$ , Figure 5bc).

451 Figure 5c breaks down the within-trial scene dissimilarity by test position for the  
 452 two age groups. Adults' scenes are clearly anchored to the initial example (right hand  
 453 facet) — shown by the dark shading in the top row indicating high similarity decreasing  
 454 from left to right for later tests — Adults' scenes are also substantially sequentially  
 455 self-similar — shown by the relatively darker shading along the diagonal compared to the  
 456 off-diagonal. In contrast, children's similarity patterns look more uniform. However, for  
 457 both adults and children, the first self-generated scene is more similar to the Initial  
 458 example than any other scene.

459 These patterns manifest in small differences in the quality of the total evidence  
 460 generated according to an information gain analysis. Adults' scenes are more informative  
 461 under a uniform prior or either the childlike or adultlike prior (see Appendix).

## 462 Experiment Discussion

463 Our constructivist analysis suggests we cannot attribute the differences in  
 464 hypothesis generation to the differences in active learning, nor can we attribute the  
 465 differences in active learning directly to differences in hypothesis generation. That is,  
 466 assuming the same generation process for the children's and then the adults' data does not  
 467 reproduce the differences in rule guesses and accuracy and assuming scene creation is  
 468 driven by distinguishing among the hypotheses manifesting the childlike or adultlike latent  
 469 prior does not explain the developmental differences in the complexity and systematicity of  
 470 the scenes creation. Rather, these data support the idea that developmental shift in  
 471 hypothesis generation and active search are manifestations a gradual tuning of the  
 472 constructivist generative mechanisms that produce novelty in cognition.

473 We now turn to a model-based analysis of the free responses and generalizations. To  
 474 foreshadow, we find both children's and adults' guesses are better accounted for by a  
 475 partially "bottom up" IDG account of hypothesis generation than by a fully top down  
 476 PCFG norm. We then find that both children's and adults' generalizations cannot be  
 477 explained by a non-symbolic similarity-based model but are well predicted by their explicit  
 478 rule guesses and by our end-to-end symbolic account, the IDG in particular.

479

## Model fitting

480 **Guesses**

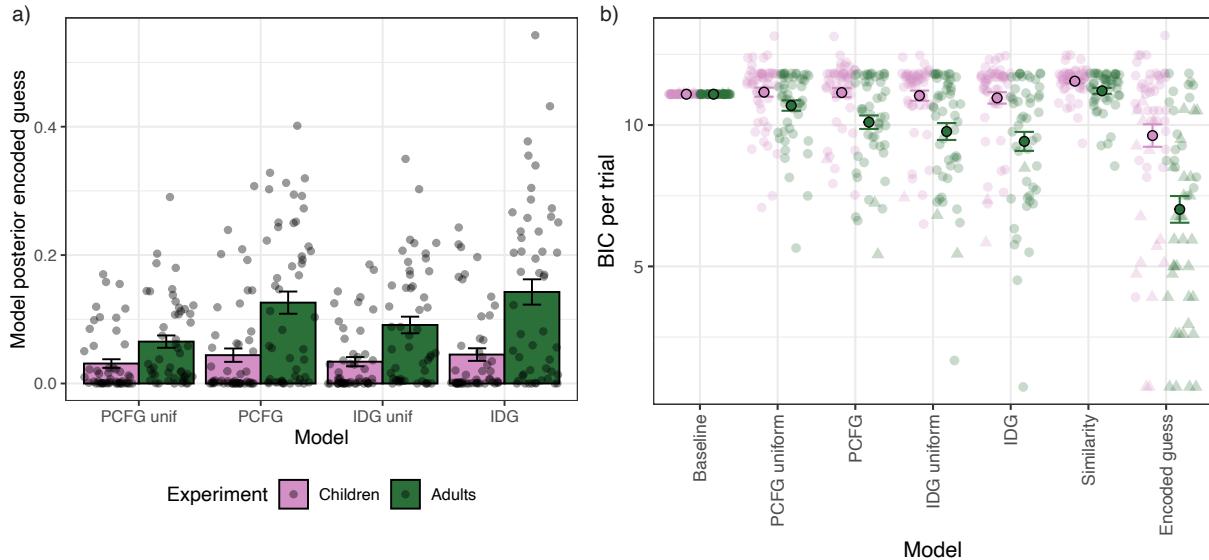
481 To evaluate our constructivist PCFG and IDG models' ability to explain  
 482 participants' free response guesses, we computed probability of each generating exactly the  
 483 participant's encoded guess conditioned on their active learning data.

484 By definition, all 87% of participants' rules that we were able to encode in our  
 485 concept grammar have nonzero support under a PCFG prior, and due to the stochasticity  
 486 we assumed in our likelihood function, all should also nonzero have posterior probability.  
 487 However, in practice it is impossible to cover an infinite space of discrete possibilities with  
 488 a finite set of samples, meaning there are a substantial number of cases in which the  
 489 complex rule produced by the participant did not appear in our PCFG or IDG samples at  
 490 all. The proportion of Children's and Adults' rules that were generated at least once in  
 491 50,000 samples by the fitted PCFG was almost identical: 68% for Children's and 67% for  
 492 Adults' guesses. The IDG fits were based on 10,000 samples generated separately for each  
 493 trial and these samples included children's rules 50% of the time and adults' rules 60% of  
 494 time using fitted weights and 49% and 62% of the time using uniform weights. Intuitively,  
 495 the slightly lower coverage is due to the practical constraint that the IDG samples had to  
 496 be computed separately for each trial limiting the number we were able to store and  
 497 evaluate and due to the constraint that the current form of our IDG algorithm produces  
 498 hypotheses with a maximum of two bound variables. The larger difference in coverage  
 499 between children's and adults guesses for the IDG is consistent with the idea that the more  
 500 complex learning scenes generated by the children resulted in a wider spread of  
 501 scene-inspired hypotheses and a concomitantly lower chance of this including the children's  
 502 explicit guess. For this reason, we visualize the posterior guess probabilities based on the  
 503 samples at the trial level (Figure 6a).

504 This reveals the IDG is the stronger hypothesis generation candidate, assigning  
 505 higher probabilities on average to the rules that participants guessed . As expected, the  
 506 variants of the PCFG and IDG with fitted production weights are better aligned with  
 507 participants' guesses than variants with uniform or mismatched weights. However, all  
 508 models produced adults' guesses with a higher probability than children's guesses.

509 **Generalizations**

510 A standard benchmark for models of concept learning is a fit with participants'  
 511 generalizations to new exemplars. Thus, we complete our analyses by comparing 18 models  
 512 ability to account for participant's generalizations. The set of models we consider allows us

**Figure 6**

a) *Encoded guess model fits:* Bars and Points show each model's per participant average ( $\pm SE$ ) posterior probability for all trials in which participants provided an encodeable guess (85%,  $N=424$ ). b) *Generalization model fits.* BIC-per-trial completed for individual level model fits to participants generalization choices (lower is better). Opaque points show model mean $\pm SE$ , faint points denote individual fits, with triangles used to mark where the model is the best fit (of all 18 tested) for that participant.

513 to test our core claims that children's and adults' induced representations are symbolic and  
 514 compositional, as opposed to statistical and similarity-driven.

515 We fit a total of 18 models to the data. All models have between 0 and 2  
 516 parameters. For each model, we fit the parameter(s) by maximizing the model's likelihood  
 517 of producing the participant data, using R's `optim` function. We compare models using the  
 518 Bayesian Information Criterion (Schwarz, 1978) to accommodate their different numbers of  
 519 fitted parameters.

520 The models we fit were:

521 **1. Baseline.** Simply assigns a likelihood of .5 to each generalization  $\in \{\text{rule}$   
 522  $\text{following}, \text{not rule following}\}$  for each of the 8 generalization probes for each of the 5  
 523 learning trials.

524 **2. Encoded Guess.** This model takes participants' free guess of the hidden rule,  
 525 coded in lambda abstraction, and uses it to generate a prediction vector  
 526  $r \in R : \{\text{rule-following}=1, \text{not rule-following}=0\}$  for each scene. These predictions are  
 527 then softmaxed using  $P(\text{choice}) = \frac{e^{r\tau}}{\sum_{r \in R} e^{r\tau}}$ , with inverse temperature parameter  
 528  $\tau \in (-\infty, \infty)$  (Luce, 1959) optimised to maximize model likelihood. Large positive  $\tau$

529 indicates a hard maximization.  $\tau \approx 0$  indicates random selection and negative  $\tau < 0$   
 530 indicates anticorrelation between model predictions and choices. For trials in which  
 531 the participant does not provide an unambiguous rule, the model assigns a .5  
 532 likelihood to each generalization choice.

533 **3. Similarity.** Inspired by Tversky's statistical and similarity based *contrast model*  
 534 of categorization (cf., Tversky, 1977), we used the inter-scene similarity between each  
 535 generalization scene and each training scene to compute the relative similarity of each  
 536 generalization case to the rule-following vs. the not rule-following training scenes.  
 537 Similarities were computed using the same procedure used in the Active Learning  
 538 section of the Results and detailed in the Appendix. We computed the mean  
 539 difference between rule-following and not-rule following similarities as a  $\Delta\text{Similarity}$   
 540 score for each participant  $\times$  trial  $\times$  item combination. Positive scores mean  
 541 generalization item has a greater feature similarity to the rule following learning  
 542 scenes than the not rule-following learning scenes. Negative scores mean the reverse.  
 543 To convert these into choice probabilities, we take the inverse logit of these scores  
 544  $r = \frac{\Delta\text{Similarity}}{\Delta\text{Similarity} + 1}$  and again fit these  $r$  values to maximize the likelihood of  
 545 participants' choices using a softmax function with inverse temperature parameter  
 546  $\tau \in (-\infty, \infty)$ . Intuitively, this model provides a non-symbolic alternative account of  
 547 generalization behavior.

548 **4-6. PCFG {uniform, off}.** These models use the marginal likelihood of each  
 549 generalization scene as rule following under the Probabilistic Context Free  
 550 Generation (PCFG) posterior as the predicted  $r$  values. “Uniform” uses the prior  
 551 with uniform production weights. “Off” uses the mismatched weights — that is,  
 552 adultlike weights for children’s generalizations and childlike weights for adults’  
 553 generalizations. As above, in each case, these prediction are fit to participants’  
 554 choices using temperature parameter  $\tau \in (-\infty, \infty)$ .

555 **7-9. IDG {uniform, off}.** These models use the marginal likelihood of each  
 556 generalization scene as rule following under the Instance Driven Generation based  
 557 posteriors with variants as with the PCFG variants and again fit with softmax  
 558 parameter  $\tau \in (-\infty, \infty)$ .

559 **10. Intercept.** This model acts a stronger baseline by allowing participants to have  
 560 an overall bias toward or against selecting generalization scenes as rule following. For  
 561 this model,  $b = 1$  if  $>50\%$  of generalizations are rule following and 0 otherwise. The

model is fit using a mixture parameter  $\lambda$  to mix this modal prediction with the baseline prediction of .5  $P(\text{choice}) = \lambda b + (1 - \lambda).5$ .

**11-18. {Other model} + Intercept.** These model variants combine the above model predictions ( $r$  = e.g. the Similarity Score or Encoded Guess prediction), with an overall generalization bias  $b$  as in 10. This involves jointly optimizing both the mixture  $\lambda$  and inverse temperature  $\tau$  parameter such that:

$$P(\text{choice}) = \lambda b + (1 - \lambda) \frac{e^{r/\tau}}{\sum_{r \in R} e^{r/\tau}}. \quad (6)$$

We fit all models to the children’s and adults’ data, and then separately to each individual participant. Individual level results are highlighted in Figure 6b. In all cases, the Intercept+PCFG and Intercept+IDG models performed no better than the pure PCFG or IDG models. Therefore we visualise individual participant fits for just models 1.–12 in Figure 6b. The full table of model fits is presented in the Appendix (Table A-3).

In line with our core hypotheses that participants’ inferences were symbolic, Encoded guess + Intercept is the best fitting model for both children’s and adults’ generalizations outperforming all the models we considered based just on only the learning data. For children, Encoded guess + Intercept has BIC 2149, improving 490 over Baseline with bias term mixture parameter  $\lambda = .26$  and choice temperature parameter  $\tau = 1.25$ . For adults, this is BIC 1776 with a larger BIC improvement of 996 over Baseline, with a  $\lambda = 0.08$  indicating lower bias and temperature  $\tau = 2.0$  indicating higher fidelity alignment with the guessed-rule’s predictions. Probing this bias, we see children undergeneralized on average, selecting just  $2.75 \pm 1.42/8$  scenes compared to adults’  $3.42 \pm 1.03/8$  — unknown to the participants, there were always 4 rule following generalization scenes. Focusing on individual fits, 16/50 children were best fit by the Baseline+Intercept model, followed by 15 by the Encoded Guess model, 9 by the Encoded Guess + Intercept model and a further 7 by Baseline. No other model best fit more than 2 children. For adults, 32/52 were best fit by Encoded guess, 6 by Baseline + Intercept, 4 by Encoded Guess + Intercept and no other model best fit more than 2 participants.

Thus, confirming our key constructivist hypothesis, there is a clear alignment between participant’s symbolic rule guesses and their generalizations. Our end-to-end constructivist sampling models (blind to the participant’s explicit guess) also received empirical support, predicting adults’ generalizations well above baselines and above the similarity-driven account which had no correlation with adults’ or children’s generalizations (see Appendix Table A-3). Thus, these results suggest participants’ inductive inferences

594 can be explained by a computational constructivism framework combined with  
 595 approximate Bayesian inference. As with the free responses, the IDG is better aligned with  
 596 participants than the PCFG, particularly when using reverse engineered rather than  
 597 uniform production weights.

## 598 Discussion

599 We explored children and adults' hypothesis generation and inductive inferences in  
 600 an interactive task where the space of possibilities and actions is open and practically  
 601 unbounded. We showed our computational constructivism framework can explain  
 602 participants successes and helps us unpack the differences between children's and adults'  
 603 behavior as consequences of differences in their generation and search mechanisms. In  
 604 particular, we found support for a partially bottom-up Instance Driven Generation account  
 605 over a fully top down Context Free (PCFG) approach, replicating (Bramley et al., 2018).  
 606 We take this to support the idea that human inductive inferences spontaneously utilize  
 607 compositional construction (Piantadosi, 2021) and further, that our introspective  
 608 descriptions pick out genuine structural features of the consequent representations. Our  
 609 formal model comparison supports these conclusions, with both end-to-end PCFG and  
 610 IDG models predicting generalizations even while blind to the free responses while a  
 611 feature-similarity model received no empirical support.

## 612 Developmental differences

613 Our analyses revealed a variety of developmental differences. Children's guesses  
 614 were more complex than adults', and consequently we could capture them with a  
 615 significantly "flatter" generation process that inherently produced a wider diversity of  
 616 hypotheses. This is broadly normative: Having been exposed to less evidence, with less  
 617 idea what conceptual compositions and fragments will be useful in understanding their  
 618 environment, children are justified in entertaining a wider set of ideas. Children were more  
 619 likely to refer to relational and positional properties in their guesses, while adults were by  
 620 most likely to make guesses that pertained to the primary object features (color and size).  
 621 This is an independently interesting finding, since relational features are structurally more  
 622 complex than primitive features, we might have predicted they would be more readily  
 623 evoked by adults. It could be that children bought in more to the scientific reasoning cover  
 624 story, treating mechanistic explanations, such as that objects must touch or be positioned  
 625 in particular ways to produce stars, as credible (Gelman, 2004). Conversely, adults may  
 626 have been more likely to expect Gricean considerations to apply, e.g. that experimenters

627 would likely set simple rules using salient but abstract features like color over perceptually  
628 ambiguous properties like position (Szollosi & Newell, 2020).

629 Children also produced more elaborate scenes during active learning than adults.  
630 One possible characterization is that children's active scene construction also used a  
631 "flatter" generative prior, resulting in more diversity of exploration approaches. Indeed,  
632 differences in active exploration are the other side of the coin of the high temperature  
633 search idea (Friston et al., 2016; Gopnik, 2020; Klahr & Dunbar, 1988; Schulz, Klenske,  
634 Bramley, & Speekenbrink, 2017). On the other hand, adults' testing behavior was more  
635 systematic, potentially reflecting a more top-down, or strategic, *control of variables*  
636 approach to gathering evidence and updating beliefs (Kuhn & Brannock, 1977). Within  
637 each trial, children's testing was more repetitive, suggesting that they made slower progress  
638 in exploring the problem space.

639 Children's guesses were also moderately less consistent with their evidence than  
640 adults'. This might be because they were less able to extract common features across  
641 learning scenes (Ruggeri & Feufel, 2015; Ruggeri & Lombrozo, 2015). However, it could  
642 also be a consequence of limited hypothesis generation. With a flatter prior and limited  
643 sampling, one has a lower chance generating a hypothesis that can explain all the evidence.  
644 Children also under-generalized, selecting only 1 or 2 of the 8 test scenes (there was  
645 actually always 4) and even when their explicit guesses predicted more should be selected.  
646 On the face of it, this reflects Wu et al.'s (Wu et al., 2017) finding that children are weaker  
647 generalizers than adults.

#### 648 Limitations and future directions

649 While this dataset provides an exceptionally rich window on developmental  
650 differences in inductive inference, some of what this task gains in open-endedness it loses in  
651 experimental control. There is residual ambiguity about the extent that differences in  
652 active learning cause differences in hypothesis generation and visa versa. Partialing this  
653 out would require controlled studies that fix the evidence and probe the hypotheses  
654 generated, or that fix the hypotheses and probe what evidence is sought. However, we have  
655 argued that constrained tasks run the risk of short-circuiting natural cognition. Learners  
656 may struggle to test hypotheses they did not conceive themselves, or to use data they have  
657 not generated to evaluate their hypotheses (Markant & Gureckis, 2014), meaning sole  
658 quantitative focus on studies that fix one or other aspect of the problem may provide a  
659 misleading perspective on active inference in the wild.

660 We also note that there are many ways we could have set up the primitives and  
661 productions of our PCFG and IDG models. Combined with non-uniform weights, this

662 makes for a dangerously expressive set of theories of cognition. We do not claim to  
663 have explored this space systematically here but that our modeling lends support to the  
664 idea that some symbolic and compositional process drives inductive inference. Identifying  
665 the computational primitives of thought may not be a realistic goal since a feature of  
666 constructivist accounts is their flexibility. Learners can grow their concept grammar over  
667 time, caching new primitives that prove useful (Piantadosi, 2021). Thus we expect different  
668 learners to take different paths in an inherently stochastic learning trajectory, limiting  
669 universal claims about representational content.

670 We assumed our scenes had directly observable features and cued these to  
671 participants in our instructions. However, a number of recent models in machine learning  
672 combine neural network methods for feature extraction with compositional engines for  
673 symbolic inference, creating hybrid systems that can learn rules and solve problems from  
674 raw inputs like natural images (cf. Nye, Solar-Lezama, Tenenbaum, & Lake, 2020; Valkov,  
675 Chaudhari, Srivastava, Sutton, & Chaudhuri, 2018). We see these approaches as having  
676 promise to bridge the gap between subsymbolic and symbolic cognitive processing.

677 Finally, as it stands, neither our PCFG or IDG are plausible process models. The  
678 PCFG is a framework for normative top-down inference in the limit of infinite sampling,  
679 and IDG is a hybrid that is less sample efficient as a brute force approach to inference  
680 in situations where a learner already has some evidence. A process account needs to  
681 explain how a learner searches the latent posterior in either framework with limited  
682 memory and computation. We see incremental adaptation of one or a few hypotheses in  
683 the light of evidence as a promising approach (Bramley, Mayrhofer, Gerstenberg, &  
684 Lagnado, 2017; Dasgupta, Schulz, & Gershman, 2017; Ullman, Goodman, & Tenenbaum,  
685 2012). A learner might use an observation to generate an initial idea akin to our IDG, but  
686 then explore permutations to this to account for the rest of the evidence. Such a process  
687 account could also help differentiate recent perspectives on the source of developmental  
688 differences. For example, it would allow measurement of whether children's search patterns  
689 are more "high temperature" than adults' (Gopnik, 2020), over and above differences in the  
690 latent search space.

## 691 Conclusions

692 We analyzed an experiment combining rich qualitative and quantitative measures of  
693 children's and adults' inductive inference. We found a number of developmental differences  
694 and demonstrated that we can make sense of these through the computational  
695 constructivism lens. Our results add empirical support and theoretical detail to recent  
696 characterizations of children as more diverse thinkers and active learners than adults, and

697 bring us closer to a computational understanding of human learning across the lifespan.

698

## References

- 699 Allen, K. R., Smith, K. A., & Tenenbaum, J. B. (2020). Rapid trial-and-error learning  
700 with simulation supports flexible tool use and physical reasoning. *Proceedings of the  
701 National Academy of Sciences*, 117(47), 29302–29310.
- 702 Anderson, J. (1990). *The adaptive character of thought*. Erlbaum.
- 703 Bramley, N. R., Dayan, P., Griffiths, T. L., & Lagnado, D. A. (2017). Formalizing  
704 Neurath's ship: Approximate algorithms for online causal learning. *Psychological  
705 Review*, 124(3), 301–338.
- 706 Bramley, N. R., Mayrhofer, R., Gerstenberg, T., & Lagnado, D. A. (2017). Causal learning  
707 from interventions and dynamics in continuous time. In *Proceedings of the 39<sup>th</sup>  
708 Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive Science  
709 Society.
- 710 Bramley, N. R., Rothe, A., Tenenbaum, J. B., Xu, F., & Gureckis, T. M. (2018).  
711 Grounding compositional hypothesis generation in specific instances. In *Proceedings  
712 of the 40<sup>th</sup> Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive  
713 Science Society.
- 714 Bruner, J. S., Jolly, A., & Sylva, K. (1976). *Play: Its role in development and evolution*.  
715 Penguin.
- 716 Carey, S. (2009). *The origin of concepts: Oxford series in cognitive development*. Oxford  
717 University Press, England.
- 718 Chen, Z., & Klahr, D. (1999). All other things being equal: Acquisition and transfer of the  
719 control of variables strategy. *Child development*, 70(5), 1098–1120.
- 720 Clarke, V., & Braun, V. (2014). Thematic analysis. In *Encyclopedia of critical psychology*  
721 (pp. 1947–1952). Springer.
- 722 Coenen, A., Nelson, J., & Gureckis, M., Todd. (2018). Asking the right questions about  
723 the psychology of human inquiry: Nine open challenges. *Psychonomic Bulletin &  
724 Review*, 26, 1548–1587.
- 725 Cook, C., Goodman, N. D., & Schulz, L. E. (2011). Where science starts: Spontaneous  
726 experiments in preschoolers' exploratory play. *Cognition*, 120(3), 341–349.
- 727 Dasgupta, I., Schulz, E., & Gershman, S. J. (2017). Where do hypotheses come from?  
728 *Cognitive psychology*, 96, 1–25.
- 729 Dennett, D. C. (1991). *Consciousness explained*. London, UK: Penguin.
- 730 Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., ... Tenenbaum, J. B.  
731 (2020). Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep  
732 bayesian program learning. *arXiv preprint arXiv:2006.08381*.
- 733 Fedyk, M., & Xu, F. (2018). The epistemology of rational constructivism. *Review of*

- 734        *Philosophy and Psychology*, 9(2), 343–362.
- 735    Feldman, J. (2000). Minimization of Boolean complexity in human concept learning.  
736              *Nature*, 407(6804), 630.
- 737    Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., Pezzulo, G., et al. (2016). Active  
738              inference and learning. *Neuroscience & Biobehavioral Reviews*, 68, 862–879.
- 739    Gelman, S. A. (2004). Psychological essentialism in children. *Trends in cognitive sciences*,  
740              8(9), 404–409.
- 741    Gershman, S. J., & Niv, Y. (2010). Learning latent structure: carving nature at its joints.  
742              *Current Opinion in Neurobiology*, 20(2), 251–256.
- 743    Ginsburg, S. (1966). *The mathematical theory of context free languages*. McGraw-Hill  
744              Book Company.
- 745    Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational  
746              analysis of rule-based concept learning. *Cognitive Science*, 32(1), 108–154.
- 747    Gopnik, A. (1996). The scientist as child. *Philosophy of science*, 63(4), 485–514.
- 748    Gopnik, A. (2020). Childhood as a solution to explore-exploit tensions. *Philosophical  
749              Transactions of the Royal Society B*, 375(1803), 20190502.
- 750    Gureckis, T. M., & Markant, D. B. (2012, September). Self-Directed Learning: A  
751              Cognitive and Computational Perspective. *Perspectives on Psychological Science*,  
752              7(5), 464–481.
- 753    Heath, C. (2004). *Zendo—Design History*. Retrieved from  
754              <http://www.koryheath.com/zendo/design-history/>
- 755    Holyoak, K. J., & Thagard, P. (1989). Analogical mapping by constraint satisfaction.  
756              *Cognitive science*, 13(3), 295–355.
- 757    Howson, C., & Urbach, P. (2006). *Scientific reasoning: the Bayesian approach*. Open  
758              Court Publishing.
- 759    Johansson, P., Hall, L., & Sikström, S. (2008). From change blindness to choice blindness.  
760              *Psychologica*, 51(2), 142–155.
- 761    Jones, A., Bramley, N. R., Gureckis, T. M., & Ruggeri, A. (in revision). Changing many  
762              things at once sometimes makes for a good experiment, and children know that.  
763              Retrieved from [psyarxiv.com/9qv5y](https://psyarxiv.com/9qv5y) doi: 10.31234/osf.io/9qv5y
- 764    Jones, M., & Love, B. C. (2011). Bayesian Fundamentalism or Enlightenment? On the  
765              explanatory status and theoretical contributions of Bayesian models of cognition. *The  
766              Behavioral and Brain Sciences*, 34(4), 169–88.
- 767    Klahr, D., & Dunbar, K. (1988). Dual space search during scientific reasoning. *Cognitive  
768              science*, 12(1), 1–48.
- 769    Klahr, D., Fay, A. L., & Dunbar, K. (1993). Heuristics for scientific experimentation: A

- 770 developmental study. *Cognitive Psychology*, 25(1), 111–146.
- 771 Klahr, D., Zimmerman, C., & Jirout, J. (2011). Educational interventions to advance  
772 children's scientific thinking. *Science*, 333(6045), 971–975.
- 773 Krippendorff, K. (2012). *Content analysis: An introduction to its methodology*. Sage.
- 774 Kruschke, J. K. (1992). Alcove: an exemplar-based connectionist model of category  
775 learning. *Psychological Review*, 99(1), 22.
- 776 Kuhn, D., & Brannock, J. (1977). Development of the isolation of variables scheme in  
777 experimental and “natural experiment” contexts. *Developmental psychology*, 13(1),  
778 9.
- 779 Lai, L., & Gershman, S. J. (2021). Policy compression: an information bottleneck in action  
780 selection.
- 781 Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building  
782 machines that learn and think like people. *Behavioral and Brain Sciences*, 40.
- 783 LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436.
- 784 Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and  
785 reversals. In *Soviet physics doklady* (Vol. 10, pp. 707–710).
- 786 Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). Sustain: a network model of category  
787 learning. *Psychological Review*, 111(2), 309.
- 788 Lucas, C. G., Bridgers, S., Griffiths, T. L., & Gopnik, A. (2014). When children are better  
789 (or at least more open-minded) learners than adults: Developmental differences in  
790 learning the forms of causal relationships. *Cognition*, 131(2), 284–299.
- 791 Luce, D. R. (1959). *Individual choice behavior*. New York: Wiley.
- 792 Markant, D. B., & Gureckis, T. M. (2014). Is it better to select or to receive? learning via  
793 active and passive hypothesis testing. *Journal of Experimental Psychology: General*,  
794 143(1), 94.
- 795 Marr, D. (1982). *Vision*. New York: Freeman & Co.
- 796 Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning.  
797 *Psychological Review*, 85(3), 207.
- 798 Michalski, R. S. (1969). On the quasi-minimal solution of the general covering problem. In  
799 (Vol. A3, pp. 125–128).
- 800 Nelson, J. D. (2005). Finding useful questions: on Bayesian diagnosticity, probability,  
801 impact, and information gain. *Psychological Review*, 112(4), 979–99.
- 802 Newell, B. R., & Shanks, D. R. (2014). Unconscious influences on decision making: A  
803 critical review. *Behavioral and brain sciences*, 37(1), 1–19.
- 804 Nye, M. I., Solar-Lezama, A., Tenenbaum, J. B., & Lake, B. M. (2020). Learning  
805 compositional rules via neural program synthesis. *arXiv preprint arXiv:2003.05562*.

- 806 Pearson, K. (1930). *The life, letters and labours of francis galton*. Cambridge University  
807 Press.
- 808 Piaget, J., & Valsiner, J. (1930). *The child's conception of physical causality*. Transaction  
809 Pub.
- 810 Piantadosi, S. T. (2021). The computational origin of representation. *Minds and  
811 Machines*, 31(1), 1–58.
- 812 Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2016). The logical primitives of  
813 thought: Empirical foundations for compositional cognitive models. *Psychological  
814 Review*, 123(4), 392.
- 815 Posner, M. I., & Keele, S. W. (1968). On the genesis of abstract ideas. *Journal of  
816 experimental psychology*, 77(3p1), 353.
- 817 Quine, W. v. O. (1969). *Word and object*. MIT press.
- 818 Ruggeri, A., & Feufel, M. (2015). How basic-level objects facilitate question-asking in a  
819 categorization task. *Frontiers in psychology*, 6, 918.
- 820 Ruggeri, A., & Lombrozo, T. (2014). Learning by asking: How children ask questions to  
821 achieve efficient search. In *36<sup>th</sup> annual meeting of the cognitive science society* (pp.  
822 1335–1340). Austin, TX: Cognitive Science Society.
- 823 Ruggeri, A., & Lombrozo, T. (2015). Children adapt their questions to achieve efficient  
824 search. *Cognition*, 143, 203–216.
- 825 Ruis, L., Andreas, J., Baroni, M., Bouchacourt, D., & Lake, B. M. (2020). A benchmark  
826 for systematic generalization in grounded language understanding. *arXiv preprint  
827 arXiv:2003.05161*.
- 828 Rule, J. S., Schulz, E., Piantadosi, S. T., & Tenenbaum, J. B. (2018). Learning list  
829 concepts through program induction. *BioRxiv*, 321505.
- 830 Rule, J. S., Tenenbaum, J. B., & Piantadosi, S. T. (2020). The child as hacker. *Trends in  
831 Cognitive Sciences*.
- 832 Schulz, E., Klenske, E. D., Bramley, N. R., & Speekenbrink, M. (2017). Strategic  
833 exploration in human adaptive control. *bioRxiv*, 110486.
- 834 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2),  
835 461–464.
- 836 Shannon, C. E. (1951). Prediction and entropy of printed english. *The Bell System  
837 Technical Journal*, 30, 50–64.
- 838 Shepard, R. N., & Chang, J.-J. (1963). Stimulus generalization in the learning of  
839 classifications. *Journal of Experimental Psychology*, 65(1), 94.
- 840 Szollosi, A., Liang, G., Konstantinidis, E., Donkin, C., & Newell, B. R. (2019).  
841 Simultaneous underweighting and overestimation of rare events: Unpacking a

- 842 paradox. *Journal of Experimental Psychology: General*, 148(12), 2207.
- 843 Szollosi, A., & Newell, B. R. (2020). People as intuitive scientists: Reconsidering statistical  
844 explanations of decision making. *Trends in Cognitive Sciences*.
- 845 Tversky, A. (1977). Features of similarity. *Psychological Review*, 84(4), 327.
- 846 Ullman, T. D., Goodman, N. D., & Tenenbaum, J. B. (2012). Theory learning as stochastic  
847 search in the language of thought. *Cognitive Development*, 27(4), 455–480.
- 848 Valkov, L., Chaudhari, D., Srivastava, A., Sutton, C., & Chaudhuri, S. (2018). Houdini:  
849 Lifelong learning as program synthesis. In *Advances in Neural Information  
850 Processing Systems* (pp. 8687–8698).
- 851 Wu, C. M., Schulz, E., Speekenbrink, M., Nelson, J. D., & Meder, B. (2017). Exploration  
852 and generalization in vast spaces. *bioRxiv*. doi: 10.1101/171371
- 853 Xu, F. (2019). Towards a rational constructivist theory of cognitive development.  
854 *Psychological Review*, 126(6), 841.

855

## Appendix A: Models

### 856 Generating context free (PCFG) model predictions

857 Here, we created a grammar (specifically a *probabilistic context free grammar* or  
 858 PCFG; Ginsburg, 1966) that can be used to produce any rule that can be expressed with  
 859 first-order logic and lambda abstraction. The grammatical primitives are detailed in  
 860 Table A-1.

**Table A-1***A Concept Grammar for the Task*

Meaning	Expression
There exists an $x_i$ such that...	$\exists(\lambda x_i : , \mathcal{X})$
For all $x_i$ ...	$\forall(\lambda x_i : , \mathcal{X})$
There exists {at least, at most, exactly} $N$ objects in $x_i$ such that...	$N_{\{<, >, =\}}(\lambda x_i : , N, \mathcal{X})$
Feature $f$ of $x_i$ has value {larger, smaller, (or) equal} to $v$	$\{<, >, \leq, \geq, =\}(x_i, v, f)$
Feature $f$ of $x_i$ is {larger, smaller, (or) equal} to feature $f$ of $x_j$	$\{<, >, \leq, \geq, =\}(x_i, x_j, f)$
Relation $r$ between $x_i$ and $x_j$ holds	$\Gamma(x_i, x_j, r)$
Booleans {and,or,not}	$\{\wedge, \vee, \neq\}(x)$
Object feature	Levels
Color	{red, green, blue}
Size	{1:small, 2:medium, 3:large}
$x$ -position	(0,8)
$y$ -position	(0,8)
Orientation	{Upright, left hand side, right hand side, strange}
Grounded	true if touching the ground
Pairwise feature	Condition
Contact	true if $x_1$ touches $x_2$
Stacked	true if $x_1$ is above and touching $x_2$ and $x_2$ is grounded
Pointing	true if $x_1$ is orientated {left/right} and $x_2$ is to $x_1$ 's {left/right}
Inside	true if $x_1$ is smaller than $x_2$ + has same $x$ and $y$ position ( $\pm 0.3$ ), false

Note that  $\{<, >, \geq, \leq\}$  comparisons only apply to numeric features (e.g., size).

861

862 There are multiple ways to implement a PCFG. Here we adopt a common approach  
 863 to set up a set of string-rewrite rules (Goodman et al., 2008). Thus, each hypothesis begins  
 life as a string containing a single *non-terminal symbol* (here,  $S$ ) that is replaced using

864 rewrite rules, or *productions*. These productions are repeatedly applied to the string,  
 865 replacing non-terminal symbols with a mixture of other non-terminal symbols and terminal  
 866 fragments of first order logic, until no non-terminal symbols remain. The productions are  
 867 so designed that the resulting string is guaranteed to be a valid grammatical expression  
 868 and all grammatical expressions have a nonzero chance of being produced. In addition, by  
 869 having the productions tie the expression to bound variables and truth statements, our  
 870 PCFG serves as an automatic concept generator. Table A-2 details the PCFG we used in  
 871 the paper.

872 We use capital letters as non-terminal symbols and each rewrite is sampled from the  
 873 available productions for a given symbol.<sup>2</sup> Because some of the productions involve  
 874 branching (e.g.,  $B \rightarrow H(B, B)$ ), the resultant string can become arbitrarily long and  
 875 complex, involving multiple boolean functions and complex relationships between bound  
 876 variables.

877 We include a variant that samples uniformly from the set of possible replacements  
 878 in each case, but we also reverse engineer a set of productions that produce exactly the  
 879 statistics the empirical samples, as described in the main text.

880 We note that it is possible, in principle, to calculate a lower bound on the prior  
 881 probability for the PCFG or IDG generating a hypothesis that a participant reported, even  
 882 if it does not occur in our sample. This can be achieved by reverse engineering the  
 883 production steps that would be needed to produce the precise encoded syntax. This is a  
 884 lower bound because it does not count semantically equivalent “phrasings” of the  
 885 hypothesis that e.g. mention features in different orders or use logically equivalent  
 886 combinations of booleans. We found that complex expressions tend to have a large number  
 887 of “phrasings”. In our sample-based approximation we implicitly treat semantically  
 888 equivalent expressions as constituting the same hypothesis but note that determining  
 889 semantic equivalence is an nontrivial aspect of constructivist inference that we do not fully  
 890 address here.

## 891 Generating instance driven (IDG) model predictions

892 We used the algorithm proposed in Bramley et al. (2018) to produce a sample of  
 893 10,000 “grounded hypotheses” for each participant and trial, splitting these evenly across  
 894 the 8 learning scenes that participant produced and tested.

---

<sup>2</sup> The grammar is not strictly context free because the bound variables ( $x_1, x_2$ , etc.) are automatically shared across contexts (e.g.  $x_1$  is evoked twice in both expressions generated in Figure 2a). We also draw feature value pairs together and conditional on the type of function they inhabit, to make our process more concise, however the same sampling is achievable in a context free way by having a separate function for every feature value, i.e. “isRed()” and sampling these directly (c.f. ?).

**Table A-2**  
*Prior Production Process*

Production	Symbol	Replacements→		
Start	$S \rightarrow$	$\exists(\lambda x_i : A, \mathcal{X})$	$\forall(\lambda x_i : A, \mathcal{X})$	$N_I(\lambda x_i : A, K, \mathcal{X})$
Bind additional	$A \rightarrow$	B	S	
Expand	$B \rightarrow$	C	$J(B, B)$	$\neg(B)$
Function	$C \rightarrow$	$= (x_i, D1)$	$I(x_i, D2)$	$= (x_i, x_j, E1)^{\mathbf{a}}$
		$I(x_i, x_j, E2)^{\mathbf{a}}$	$\Gamma(x_i, x_j, E3)^{\mathbf{a}}$	
Feature/value (numeric only)	$D1 \rightarrow$	value,	feature	
	$D2 \rightarrow$	value,	feature	
Feature (numeric only)	$E1 \rightarrow$	feature		
	$E2 \rightarrow$	feature		
(relational)	$E3 \rightarrow$	feature		
Boolean	$J \rightarrow$	$\wedge$	$\vee$	...
Inequality	$I \rightarrow$	$\leq$	$\geq$	$>$
		$<$		
Number	$K \rightarrow$	$n \in \{1, 2, 3, 4, 5, 6\}$		

Note: Context-sensitive aspects of the grammar: <sup>a</sup>Bound variable(s) sampled uniformly without replacement from set; expressions requiring multiple variables censored if only one.

895 To generate hypotheses as candidates for the hidden rule, the model uses the  
 896 following procedure with probabilities either set to uniform or drawn from the PCFG-fitted  
 897 productions for adults or for children (Figure 3gh) and denoted with square brackets:

898 1. **Observe.** either:

- 899 (a) With probability  $[A \rightarrow B]$ : Sample a cone from the observation, then sample  
 900 one of its features  $f$  with probability  $[G \rightarrow f]$  — e.g.,  $\{\#1\}$ :<sup>3</sup> “medium, size” or  
 901  $\{\#3\}$ : “red, color”.  
 902 (b) With probability  $[A \rightarrow \text{Start}]$ : Sample two cones uniformly without replacement  
 903 from the observation, and sample any shared or pairwise feature — e.g.,  
 904  $\{\#1, \#2\}$ : “size”, or “contact”

905 2. **Functionize.** Bind a variable for each sampled cone in Step 1 and sample a true  
 906 (in)equality statement relating the variable(s) and feature:

- 907 (a) For a statement involving an unordered feature there is only one possibility —  
 908 e.g.,  $\{\#3\}$ : “ $= (x_1, \text{red}, \text{color})$ ”, or for  $\{\#1, \#2\}$ : “ $= (x_1, x_2, \text{color})$ ”  
 909 (b) For a single cone and an ordered feature, this could also be a nonstrict  
 910 inequality ( $\geq$  or  $\leq$ ). We assume a learner only samples an inequality if it

<sup>3</sup> Numbers prepended with # refer to the labels on the cones in the example observation in Figure 2b.

expands the number of cones picked out from the scene relative to an equality — e.g., in Figure 2b in the main text, there is also a large cone  $\{\#1\}$  so either  $\geq(x_1, \text{medium}, \text{size})$  or  $=x_1, \text{medium}, \text{size}$ ) might be selected with uniform probability.

- (c) For two cones and an ordered feature, either strict or non-strict inequalities could be sampled if the cones differ on the sampled feature, equivalently either equality or non-strict inequality could be selected if the cones do not differ on that dimension — e.g.,  $\{\#1, \#2\} > (x_1, x_2, \text{size})$ , or  $\{\#3, \#4\} \geq (x_1, x_2, \text{size})$ . In each case, the production weights from Figure 3g&h for the relevant completions are normalized and used to select the option.

3. **Extend.** With probability  $\frac{[B \rightarrow D]}{[B \rightarrow D] + [B \rightarrow C(B, B)]}$  go to Step 4, otherwise sample a conjunction with probability  $[C(B, B) \rightarrow \text{And}]$  or a disjunction with probability  $[C(B, B) \rightarrow \text{Or}]$  and repeat. For statements with two bound variables, Step 3 is performed for  $x_1$ , then again for  $x_2$ :

- (a) **Conjunction.** A cone is sampled from the subset picked out by the statement thus far and one of its features sampled with probability  $[G \rightarrow f]$  — e.g.,  $\{\#1\} \wedge (=x_1, \text{green}, \text{color}), \geq(x_1, \text{medium}, \text{size})$ ). Again, inequalities are sampleable only if they increase the true set size relative to equality — e.g., “ $\wedge(\leq(x_1, 3, \text{xposition}), \geq(x_1, \text{medium}, \text{size}))$ ”, which picks out more objects than “ $\wedge(=x_1, 3, \text{xposition}), \geq(x_1, \text{medium}, \text{size})$ ”.
- (b) **Disjunction.** An additional feature-value pair is selected uniformly from *either* unselected values of the current feature, *or* from a different feature — e.g.,  $\vee(=x_1, \text{color}, \text{red}), (=x_1, \text{color}, \text{blue})$  or  $\vee(=x_1, \text{color}, \text{blue}), \geq(x_1, \text{size}, 2)$ ). This step is skipped if the statement is already true of all the cones in the scene.<sup>4</sup>

4. **Flip.** If the inspiration scene is not rule following wrap the expression in a  $\neg()$ .

5. **Quantify.** Given the contained statement, select true quantifier(s):

- (a) For statements involving a single bound variable (i.e., those inspired by a single cone in Step 1) the possible quantifiers simply depend on the number of the cones in the scene for which the statement holds. If the statement is true of all cones in the scene Quantifier is selected using probabilities [Start→] combined with  $[L \rightarrow]$  where appropriate. If it is true of only a subset of the cones then

---

<sup>4</sup> We rounded positional features to one decimal place in evaluating rules to allow for perceptual uncertainty.

942        $\forall(\lambda x_i : A, \mathcal{X})$  is censored and the probabilities re-normalized.  $K$  is set to match  
 943       number of cones for which the statement is true.

- 944       (b) Statements involving two bound variables in lambda calculus have two nested  
 945       quantifier statements each selected as in (a). The inner statement quantifying  $x_2$   
 946       is selected first based on truth value of the expression while taking  $x_1$  to refer to  
 947       the cone observed in ‘1.’. The truth of the selected inner quantified statement is  
 948       then assessed for all cones to select the outer quantifier — e.g.,  $\{\#3, \#4\}$   
 949       “ $\wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size}))$ ” might become  
 950       “ $\forall(\lambda x_1 : \exists(\lambda x_2 : \wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size})), \mathcal{X}), \mathcal{X})$ ”. The inner  
 951       quantifier  $\exists$  is selected (three of the four cones are green  $\{\#1, \#2, \#4\}$ ), and  
 952       the outer quantifier  $\forall$  is selected (all cones are less than or equal in size to a  
 953       green cone).

954       Note that a procedure like the one laid out above is, in principle, capable of  
 955       generating any rule generated by the PCFG in Figure

956       One way to think of the IDG procedure is as a partial inversion of a PCFG. As  
 957       illustrated by the blue text in the examples in Figure 2b in the main text. While the  
 958       PCFG starts at the outside and works inward, the IDG starts from the central content and  
 959       works outward out to a quantified statement, ensuring at each step that this final  
 960       statement is true of the scene.

## 961 Full generalization model fits

962       As described in main text, we fit 18 model variants to participant’s data. All models  
 963       have between 0 and 2 parameters. For each model, we fit the parameter(s) by maximizing  
 964       the model’s likelihood of producing the participant data, using R’s `optim` function. We  
 965       compare models using the Bayesian Information Criterion (Schwarz, 1978) to accommodate  
 966       their different numbers of fitted parameters.<sup>5</sup> Full results are in Table A-3.

## 967 Appendix B: Free response coding

968       To analyze the free responses, we first had two coders go through all responses and  
 969       categorize them as either:

---

<sup>5</sup> On one perspective, our derivation of the child-like and adult-like productions constitutes fitting an additional 39 parameters ( $m - 1$  for each production step), so evoking an additional BIC parameter penalty of  $39 \times \log(3940) = 323$  for PCFG over PCFG Uniform and similarly for the IDG. If we were to apply this penalty, the uniform weighted variants would be clearly preferred under the BIC criterion at the aggregate level. It is less clear how to apply this penalty at the individual level. We chose to include the fitted versions alongside the uniform versions here without penalty as demonstrations of the differences that arise from different generation probabilities.

**Table A-3**  
*Models of Participants' Generalizations*

Model	Group	log(Likelihood)	BIC	$\lambda$	$\tau$	N	Accuracy
1. Baseline	children	-1319.75	2639.50			7	50%
2. Encoded Guess	children	-1143.69	2294.92	0.98	15	62%	
3. Similarity	children	-1316.44	2640.42	-0.50	0	41%	
4. PCFG Uniform	children	-1319.75	2647.05	-0.01	0	60%	
5. PCFG Off	children	-1318.85	2645.26	0.09	0	65%	
6. PCFG	children	-1319.57	2646.69	0.04	1	63%	
7. IDG Uniform	children	-1299.72	2607.00	0.55	2	66%	
8. IDG Off	children	-1304.92	2617.39	0.45	1	<b>70%</b>	
9. IDG	children	-1308.52	2624.59	0.39	2	68%	
10. Intercept	children	-1218.96	2445.47	0.32		<b>16</b>	50%
11. <b>Encoded Guess + Intercept</b>	<b>children</b>	<b>-1067.18</b>	<b>2149.47</b>	0.26	1.24	9	
12. Similarity + Intercept	children	-1214.71	2444.52	0.32	-0.77	1	
13. PCFG Uniform + Intercept	children	-1210.30	2435.70	0.35	0.43	0	
14. PCFG Off + Intercept	children	-1207.64	2430.39	0.34	0.48	0	
15. PCFG + Intercept	children	-1208.74	2432.59	0.35	0.46	0	
16. IDG Uniform + Intercept	children	-1195.19	2405.48	0.32	0.83	0	
17. IDG Off + Intercept	children	-1193.01	2401.12	0.34	0.83	0	
18. IDG + Intercept	children	-1194.19	2403.49	0.34	0.82	0	
1. Baseline	adults	-1386.29	2772.59			2	50%
2. Encoded Guess	adults	-893.49	1794.58	1.78		<b>32</b>	70%
3. Similarity	adults	-1359.05	2725.70	-1.38	0	36%	
4. PCFG Uniform	adults	-1333.95	2675.50	0.69	0	62%	
5. PCFG Off	adults	-1293.60	2594.79	0.94	1	66%	
6. PCFG	adults	-1267.89	2543.38	1.06	2	69%	
7. IDG Uniform	adults	-1229.69	2466.97	1.50	2	69%	
8. IDG Off	adults	-1208.11	2423.83	1.52	0	73%	
9. IDG	adults	-1185.64	2378.89	1.62	1	<b>74%</b>	
10. Intercept	adults	-1364.90	2737.40	0.15		6	50%
11. <b>Encoded Guess + Intercept</b>	<b>adults</b>	<b>-880.59</b>	<b>1776.38</b>	0.08	2.01	4	
12. Similarity + Intercept	adults	-1337.55	2690.30	0.14	-1.63	0	
13. PCFG Uniform + Intercept	adults	-1268.87	2552.93	0.26	1.35	0	
14. PCFG Off + Intercept	adults	-1226.61	2468.42	0.25	1.60	0	
15. PCFG + Intercept	adults	-1203.66	2422.53	0.24	1.69	0	
16. IDG Uniform + Intercept	adults	-1179.02	2373.24	0.20	2.13	0	
17. IDG Off + Intercept	adults	-1147.46	2310.13	0.22	2.26	0	
18. IDG + Intercept	adults	-1131.92	2279.04	0.20	2.28	0	

NB: Accuracy column shows performance of the requisite model across 100 simulated runs through the task using participants active learning data and  $\tau$  set to 100 (essentially hard maximizing over the model's predictions). The +Intercept models perform strictly worse due to their bias so are not included in this column.

- 970 1. Correct: The subject gives exactly the correct rule or something logically equivalent
- 971 2. Overcomplicated: The subject gives a rule that over-specifies the criteria needed to  
972 produce stars relative to the ground truth. This means the rule they give is logically  
973 sufficient but not necessary. For example, stipulating that "there must be a small  
974 red" is overcomplicated if the true rule is "there must be a red" because a scene could  
975 contain a medium or large red and emit stars.
- 976 3. Overliberal: The opposite of overcomplicated. The subject gives a rule that  
977 under-specifies what must happen for the scene to produce stars. For example,  
978 stipulating that "there must be a blue" if the true rule is that "exactly one is blue".

979        This is logically necessary but not sufficient because a scene could contain blue  
 980        objects but not produce stars because there is not exactly one of them.

- 981        4. Different: The subject gives a rule that is intelligible but different from the ground  
 982        truth in that it is neither necessary or sufficient for determining whether a scene will  
 983        produce stars.
- 984        5. Vague or multiple. Nuisance category.
- 985        6. No rule. The subject says they cannot think of a rule.

986        We were able to encode 205/238 (86%) of the children's responses and (219/250)  
 987        87% for adults as correct, overcomplicated, overliberal or different. Table A-4 shows the  
 988        complete confusion matrix. The two coders agreed 85% of the time, resulting in a Cohen's  
 989        Kappa of .77 indicating a good level of agreement (Krippendorff, 2012).

**Table A-4**  
*Agreement Matrix for Independent Coders' Free Response Classifications*

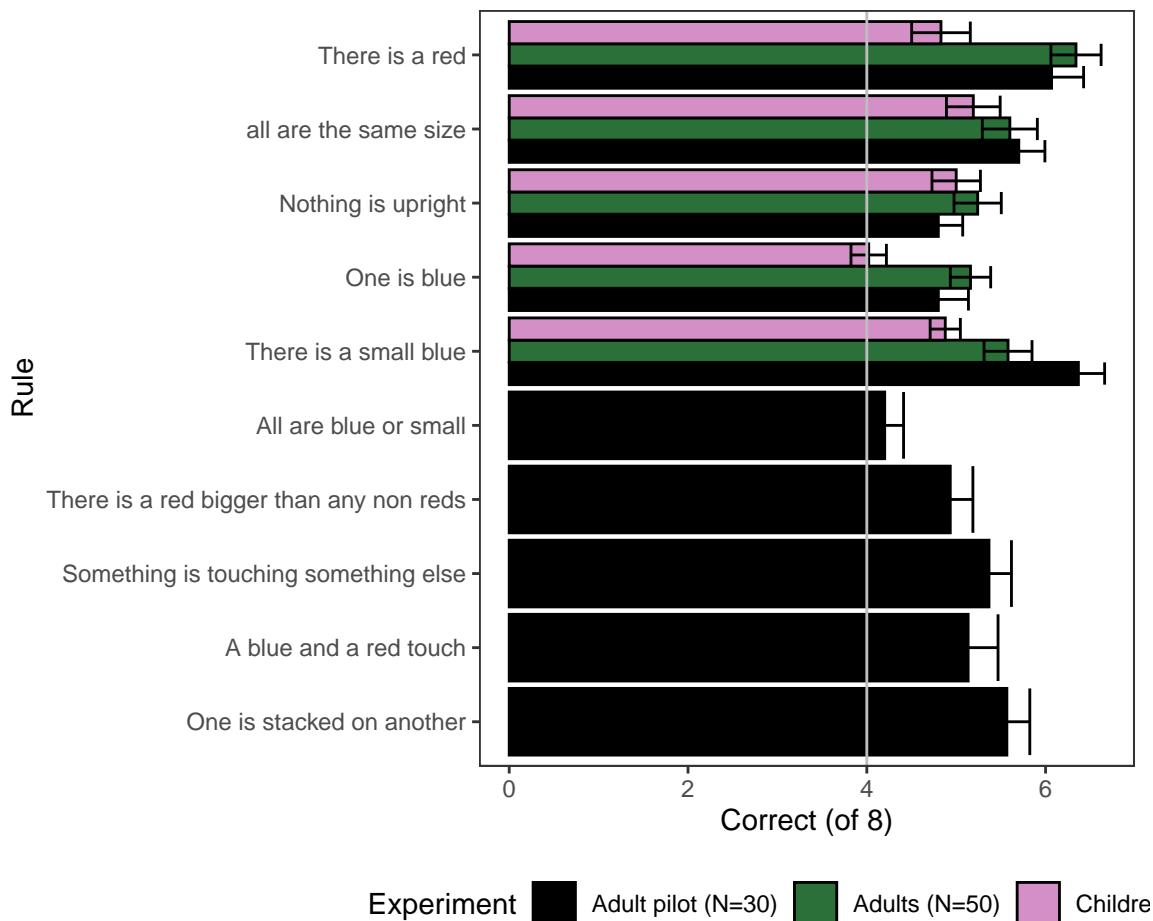
	correct	overliberal	overspecific	different	vague	no rule	multiple
correct	<b>93</b>	1	5	0	0	0	0
overliberal	5	<b>13</b>	1	8	0	1	0
overspecific	1	2	<b>42</b>	12	0	0	0
different	0	5	3	<b>224</b>	15	3	0
vague	0	1	2	3	<b>11</b>	6	0
no rule	0	0	0	0	0	<b>31</b>	0
multiple	0	1	0	2	0	0	<b>0</b>

990        We then had one coder familiar with the grammar go through each free response  
 991        that was not assigned vague or no rule, and encode it as a function in our grammar. The  
 992        second coder then blind spot checked 15% of these rules (64) and agreed in 95% of cases  
 993        61/64. The 6 cases of disagreement were discussed and resolved. In 5/6 cases, this was in  
 994        favor of the primary coder. The full set of free text responses along with the requisite  
 995        classification, encoded rules are available in the [Online Repository](#).

## 996        Appendix C: Comparison with Bramley et al (2018)

997        Finally, for interest and to demonstrate replication of our core results. We provide a  
 998        direct comparison between the generalization accuracies in the current sample of children  
 999        and adults and those in the sample of 30 adults modelled in (Bramley et al., 2018).  
 1000      Bramley et al (2018) included 10 ground truth concepts, and the current paper uses just

1001 the first five of these. Figure A-1 shows these accuracy patterns side by side revealing the  
 1002 adults in the current experiment performed approximately as well as those in the original  
 1003 conference paper.



**Figure A-1**

*Generalization accuracy by number of objects per test scene comparing with 10 rule adult pilot from Bramley et al. (2018).*

1004

#### Appendix D: Scene similarity measurement

1005 To establish the overall similarity between two scenes, we need to map the objects  
 1006 in a given scene to the objects in another scene (for example between the scenes in  
 1007 FigureA-2 a and b) and establish a reasonable cost for the differences between objects  
 1008 across dimensions. We also need a procedure for cases where there are objects in one scene  
 1009 that have no analogue in the other. We approach the calculation of similarity via the  
 1010 principle of minimum edit distance (Levenshtein, 1966). This means summing up the  
 1011 elementary operations required to convert scene (a) into scene (b) or visa versa. We assume

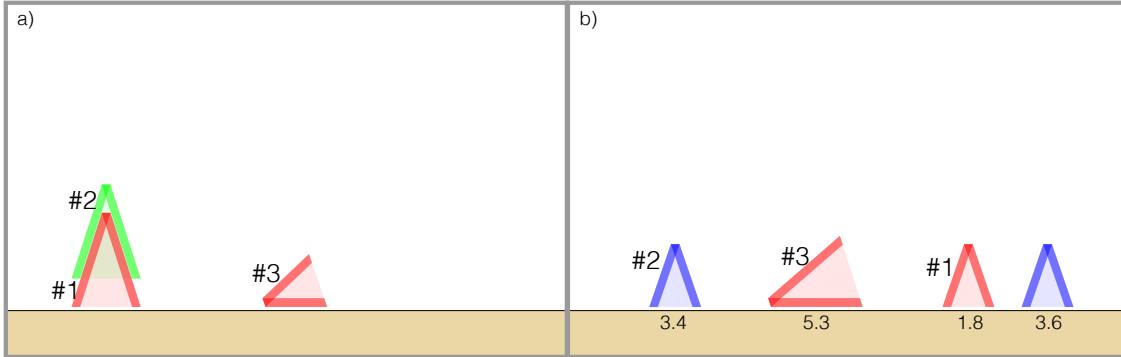
1012 objects can be adjusted in one dimension at a time (i.e. moving them on the  $x$  axis,  
 1013 rotating them, or changing their color, and so on.

1014 Before focusing on how to map the objects between the scenes we must decide how  
 1015 to measure the adjustment distance for a particular object in scene a to its supposed  
 1016 analogue in scene b. As a simple way to combine the edit costs across dimensions we first  
 1017 Z-score each dimension, such that the average distance between any two values across all  
 1018 objects and all scenes and dimensions is 1. We then take the L1-norm (or city block  
 1019 distance) as the cost for converting an object in scene (a) to an object in scene (b), or visa  
 1020 versa. Note this is sensitive the size of the adjustment, penalizing larger changes in  
 1021 position, orientation or size more severely than smaller changes, while changes in color are  
 1022 all considered equally large since color is taken as categorical. Note also that for  
 1023 orientation differences we also always assume the shortest distance around the circle.

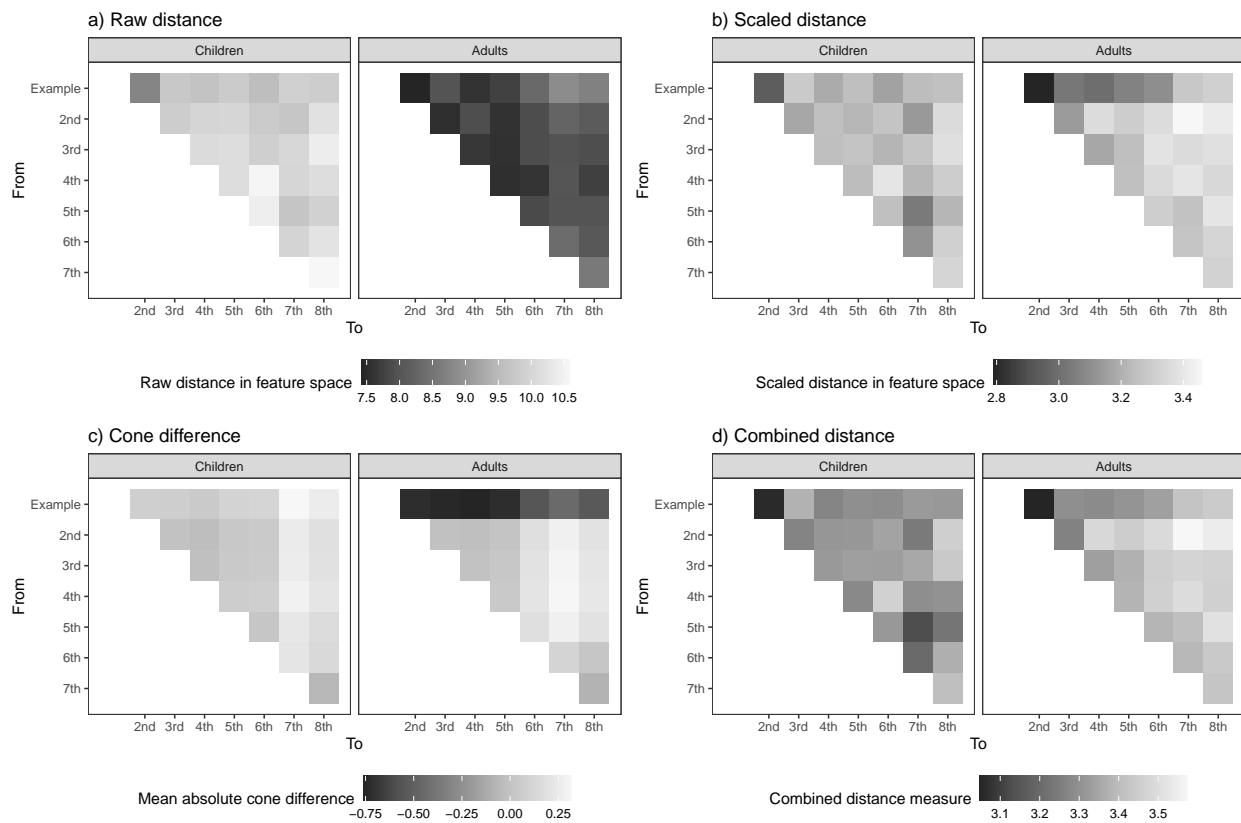
1024 If scene (a) has an object that does not exist in scene (b) we assume a default  
 1025 adjustment penalty equal to the average divergence between two objects across all  
 1026 comparisons (3.57 in the current dataset). We do the same for any object that exists in (a)  
 1027 but not (b).

1028 Calculating the overall similarity between two scenes involves solving a mapping  
 1029 problem of identifying which objects in scene (a) are “the same” as those in scene (b). We  
 1030 resolve this “charitably”, by searching exhaustively for the mapping of objects in scene (a)  
 1031 to scene (b) that minimizes the total edit distance. Having selected this mapping, and  
 1032 computed the final edit distance including any costs for additional or removed objects, we  
 1033 divide by the number shared cones, so as to avoid the dissimilarities increasing with the  
 1034 number of objects involved.

1035 Figure A-3 computes the inter-scene similarity components that go into Figure 6c in  
 1036 the main text. Summing up the edit distances across all objects, children’s scenes seem  
 1037 much more diverse than adults (Figure A-3a). However this is primarily due to their  
 1038 containing a greater average number of objects. Scaling the edit distance by the number of  
 1039 objects in the target scene gives a more balanced perspective (Figure A-3b) but does not  
 1040 account for the fact that the compared scene may contain more or fewer objects in total.  
 1041 Figure A-3c visualizes just the object difference showing that children’s scenes contain  
 1042 roughly as many objects on average as the initial example while adults’ scenes contain  
 1043 around 0.75 fewer objects than are present in the initial example (dark shading in top row).  
 1044 Thus, we opted to combine b and c by weighting the unsigned cone difference by the mean  
 1045 inter-object distance across all comparisons to give our combined distance measure  
 1046 (Figure A-3d and Figure 6c in the main text).

**Figure A-2**

Three example scenes. Objects indices link the most similar set of objects in b to those in a. Numbers below indicate the edit distance for each object (i.e. the sum of scaled dimension adjustments). Intuitively scene a) is more similar to scene b) than to scene c) and this is reflected in the similarity scores.

**Figure A-3**

a) The average minimum edit distance summed up across shared objects. b) Rescaling a by dividing by the number of objects. c) The penalty for additional or omitted objects. d) Combined distance as in main text.

**Appendix E: Information gain analysis of active learning data**

Children's and adults' scene generation patterns manifest in small differences in the

quality of the total evidence generated according to an information gain analysis. For

example, using the unweighted PCFG sample, prior entropy is 13.31 bits and children's

evidence produces an information gain (reduction in uncertainty) of  $6.86 \pm 0.55$  bits while

adults data allows for marginally higher information  $7.04 \pm 0.44$  bits

$t(102) = -1.8, p = 0.068$ . Relative to the fitted PCFG priors, the difference in information

gains is rather larger, with children's scenes leading to information gain at  $6.92 \pm 0.70$  bits

(prior entropy 12.94), and adults' at  $7.50 \pm 0.66$  (prior entropy 12.65)

$t(102) = 4.4, p < .0001$ . Under the mismatched priors — that is, taking the adultlike

PCFG prior for children and childlike PCFG prior for adults — children's tests look

slightly more informative than under their own prior, generating  $7.14 \pm 0.72$  bits, and adults'

tests slightly less informative than under their own prior  $7.21 \pm 0.61$  bits, eliminating the

statistical difference  $t(102) = 0.5, p = 0.62$ . On the face of it, this is evidence against the

idea that children's more elaborate hypothesis generation and concomitantly flatter latent

prior is driving them rationally toward more elaborate testing patterns. However, we see

this information-theoretic analyses as limited in what reveals. This is because is predicated

on an implausibly complete representation of uncertainty, e.g. using a large sample of prior

hypotheses, while we might expect constructivist search behavior to be driven by more

focal testing of a smaller number of possibilities. Nevertheless, we present these information

scores as norms for completeness and comparison with other active learning tasks.