

**Active inductive inference in children and adults: A constructivist perspective**

Neil R. Bramley\*

Department of Psychology, University of Edinburgh, Scotland

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 explicit 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 causal rules through active testing. Children were more elaborate in their testing behavior and generated substantially more complex guesses about the hidden rules. We take a ‘computational constructivist’ perspective to explaining these patterns, arguing that these inferences are driven by a combination of thinking (generating and modifying symbolic concepts) and exploring (discovering and investigating patterns in the physical world). We show how this framework and rich new dataset speak to questions about developmental differences in hypothesis generation, active learning and inductive generalization. In particular, we find children’s learning is driven by less fine-tuned construction mechanisms than adults’, resulting in a greater diversity of ideas but less reliable discovery of simple explanations.

## Active inductive inference in children and adults: A constructivist perspective

*“We think we understand the rules when we become adults but what we really experience is a narrowing of the imagination.”* — David Lynch

1       A central question in the study of both human development and reasoning is how  
2 learners come up with the ideas and hypotheses they use to explain the world around  
3 them. Children excel at forming new categories, concepts, and causal theories (Carey,  
4 2009) and by maturity, this coalesces into a capacity for intelligent thought characterized  
5 by its domain generality and occasional moments of insight and innovation. Constructivism  
6 is an influential perspective in developmental psychology (Carey, 2009; Piaget, 2013; Xu,  
7 2019) and philosophy of science (Fedyk & Xu, 2018; Phillips, 1995; Quine, 1969) that  
8 posits learners actively construct new ideas through a mixture of thinking—recombining  
9 and modifying ideas—and play—exploring and discovering patterns in the world (Bruner,  
10 Jolly, & Sylva, 1976; Piaget & Valsiner, 1930; Xu, 2019). While the tenets and promise of  
11 constructivist accounts are appealing, it has historically lacked the formalization needed to  
12 distinguish it from alternative accounts of learning, limiting its testable predictions or  
13 detailed insights into cognition. We draw on recent methodological advances to formalize  
14 key aspects of constructivism and use these to analyze children and adults’ behavior in an  
15 open-ended inductive learning task. We show that a virtue of the constructivist account is  
16 that it captures the wide range of ideas and testing behaviors we observe, particularly in  
17 children. We use our account to examine developmental differences in hypothesis  
18 generation and active learning. To foreshadow, we show children’s hypothesis generation  
19 and active learning are driven by less fine-tuned construction mechanisms than adults’,  
20 resulting in a greater diversity of ideas but less reliable discovery of simple explanations  
21 and less systematic coverage of the data space.

### 22 Concept learning

23       Classic work in experimental psychology suggests symbol manipulation is required  
24 for humanlike reasoning and problem solving (Bruner, Goodnow, & Austin, 1956;  
25 Johnson-Laird, 1983; Wason, 1968). However, classic symbolic accounts struggled to  
26 explain how discrete representations could be learned or effectively applied to reasoning  
27 under uncertainty (Oaksford & Chater, 2007; Posner & Keele, 1968). Meanwhile, statistical  
28 accounts of concept learning have flourished by treating concepts as driven by “family  
29 resemblance” within a feature space—for instance, centered around a prototypical example  
30 or set of exemplars (Kruschke, 1992; Love, Medin, & Gureckis, 2004; Medin & Schaffer,  
31 1978; Shepard & Chang, 1963). Such accounts help explain how people assign category  
32 membership fuzzily, and generalize effectively to novel stimuli (Shepard, 1987) but lack a

<sup>33</sup> core representation capable of capturing how people construct conceptual novelty  
<sup>34</sup> (Komatsu, 1992).

<sup>35</sup> Bayesian approaches have also played a major role in study of concept learning,  
<sup>36</sup> providing a principled way of modeling probabilistic inference over both sub-symbolic and  
<sup>37</sup> symbolic hypothesis spaces (Howson & Urbach, 2006). On the symbolic side this includes  
<sup>38</sup> inferences about particular causal structures (Bramley, Lagnado, & Speekenbrink, 2015;  
<sup>39</sup> Coenen, Rehder, & Gureckis, 2015; Gopnik et al., 2004; Steyvers, Tenenbaum,  
<sup>40</sup> Wagenmakers, & Blum, 2003) as well as more general causal theories (Goodman, Ullman,  
<sup>41</sup> & Tenenbaum, 2011; Griffiths & Tenenbaum, 2009; Kemp & Tenenbaum, 2009; Lucas &  
<sup>42</sup> Griffiths, 2010). Alongside Bayesian analyses, information theory has also featured  
<sup>43</sup> frequently as a metric of idealized evidence acquisition (Gureckis & Markant, 2012),  
<sup>44</sup> including choice of interventions and experiments that reveal causal structure (Bramley,  
<sup>45</sup> Dayan, Griffiths, & Lagnado, 2017; Bramley et al., 2015; Coenen et al., 2015; Steyvers et  
<sup>46</sup> al., 2003). However, since idealized Bayesian and information theoretic accounts describe  
<sup>47</sup> learning within a predefined hypothesis space, they do not directly explain how a learner  
<sup>48</sup> explores or generates possibilities within an infinite latent space. That is, probabilistic  
<sup>49</sup> accounts of induction on are generally cast at Marr's computational level (Marr, 1982),  
<sup>50</sup> showing people behave roughly *as if* they consider and average exhaustively over what is  
<sup>51</sup> really an unbounded space of possible concepts. Thus, while these accounts provide a  
<sup>52</sup> jumping off point for rational analysis of cognition, we should take their limitations  
<sup>53</sup> seriously when seeking to reverse engineer humanlike inductive inference (Simon, 2013;  
<sup>54</sup> Van Rooij, Blokpoel, Kwisthout, & Wareham, 2019).

<sup>55</sup> The goal of this paper is to examine children's and adults' inductive learning in a  
<sup>56</sup> rich open-ended task where the space of potential hypotheses and behaviors is effectively  
<sup>57</sup> unbounded. In doing this, we will treat constructivism as a form of rational process  
<sup>58</sup> framework (Lieder & Griffiths, 2020), capturing how people are shaped by Bayesian and  
<sup>59</sup> information-theoretic norms but also why they diverge from and fall short of them outside  
<sup>60</sup> of constrained scenarios. To do this, we focus on recent work in cognitive science that has  
<sup>61</sup> attempted to marry symbolic and statistical perspectives. This work characterizes  
<sup>62</sup> computational principles driving both human development and intelligence as resting on a  
<sup>63</sup> capacity to flexibly generate, adapt, combine and re-purpose symbolic representations  
<sup>64</sup> when learning and reasoning, but crucially to do so in ways that approximate probabilistic  
<sup>65</sup> principles of inference under uncertainty (Bramley, Dayan, et al., 2017; Goodman,  
<sup>66</sup> Tenenbaum, Feldman, & Griffiths, 2008; Piantadosi, 2021; Piantadosi, Tenenbaum, &  
<sup>67</sup> Goodman, 2016).

## 68 Constructivism

69       Fundamentally, we take the constructivist account to depart from  
70 computational-level Bayesian accounts because it presumes representational  
71 *incompleteness*, and consequently *stochasticity* and *path dependence* in a given individual's  
72 learning trajectory. By this, we mean that the constructivist learner has not, and normally  
73 could not, consider and weigh all the possibilities in play when learning. Instead, they  
74 must have some mechanism for generating and comparing finite numbers of discrete  
75 possibilities (Sanborn & Chater, 2016; Stewart, Chater, & Brown, 2006). Eponymously, the  
76 construction mechanism needs to be capable of recursive *construction*: composing and  
77 recomposing symbolic elements so as to achieve the systematicity and productivity  
78 required for a finite system to cover an infinite space of ideas (Piantadosi & Jacobs, 2016).  
79 In this way, constructivist views treat algorithmic-level cognition as necessarily symbolic  
80 and at least somewhat language-like (Fodor, 1975) in its ability to make "infinite use of  
81 finite means" (von Humboldt, 1863/1988).

82       For example, a constructivist learner might stochastically combine elements from an  
83 underlying concept grammar to produce new ideas that can be tested against evidence.  
84 Alternatively, they might use their grammar to describe patterns in evidence or to adapt a  
85 previous hypotheses to fit some new evidence (Bonawitz, Denison, Gopnik, & Griffiths,  
86 2014; Lewis, Perez, & Tenenbaum, 2014; Nosofsky & Palmeri, 1998; Nosofsky, Palmeri, &  
87 McKinley, 1994). Outside of narrow experimental settings, this modal incompleteness  
88 seems completely normal. A simple illustration is the gap between ease of evaluation versus  
89 generation of hypotheses (Gettys & Fisher, 1979). We can typically generate fewer  
90 explanations on the fly—i.e., reasons why our car won't start—than we would endorse if a  
91 list was presented to us. We would likely come up with more as we looked under the hood  
92 than we would sat in the car thinking. Inference about any area of active scientific inquiry,  
93 like that reported in this journal, typically involve an enormous latent space of potential  
94 explanatory theories only a fraction of which have ever been articulated or tested and  
95 many of which were discovered only serendipitously (Shackle, 2015). It is generally  
96 accepted that the ground truth is unlikely to be among the set of theories already on the  
97 table (Box, 1976) and that challenging results are as likely to lead to theory modification  
98 as complete abandonment (Lakatos, 1976).

99       The constructivist perspective thus departs from a Bayesian analysis by emphasizing  
100 that induction is as much about constructing candidate possibilities, as optimizing within a  
101 set of candidates. This reframing demystifies a number of behavioral patterns that look like  
102 biases from the computational-level perspective. These include *anchoring*, *order effects*,  
103 *probability matching* and *confirmation bias*. For example, *Anchoring* effects in estimation

104 can be explained as resulting from limited local adjustment from a salient starting point  
105 (Griffiths, Lieder, & Goodman, 2015; Lieder, Griffiths, Huys, & Goodman, 2018). *Order*  
106 *effects*, where the sequence of evidence encountered affects the final belief, are pervasive in  
107 human learning. If new hypotheses are arrived at through a limited local search starting  
108 from a previous hypothesis then we should expect patterns of path dependence and  
109 auto-correlation between a single learner's hypotheses over time (Bramley, Dayan, et al.,  
110 2017; Dasgupta, Schulz, & Gershman, 2016; Fränken, Theodoropoulos, & Bramley, 2022;  
111 Thaker, Tenenbaum, & Gershman, 2017; Zhao, Lucas, & Bramley, 2022). *Probability*  
112 *matching* is also natural under a constructivist perspective. In experiments, participants  
113 often choose options in proportion to their probability of being correct or optimal rather  
114 than reliably selecting the best action, as we might expect if they had the full posterior to  
115 hand (Shanks, Tunney, & McCarthy, 2002). However, it can be shown that rather than  
116 being a choice pathology, probability matching may be better seen as a *best case* scenario  
117 for a learner limited to using the endpoint of a local search as their guess (Bramley,  
118 Dayan, et al., 2017). It has been argued that in a variety of plausible everyday settings, a  
119 single-sample-based decision can be the appropriate computation-accuracy tradeoff for a  
120 resource-limited learner (Vul, Goodman, Griffiths, & Tenenbaum, 2009). *Confirmation bias*  
121 is also pervasive in human reasoning and active learning (Klayman & Ha, 1989) and hard  
122 to explain in purely Bayesian terms. Wason (1960) famously asked participants to test and  
123 identify a hidden rule and initially simply told them that the sequence 2–4–6 followed the  
124 rule. The intended true rule was simply “ascending numbers” but participants frequently  
125 guessed more complex rules such as “numbers increasing by two”. Analysis of participants’  
126 tests revealed that they frequently generated tests that would be rule-following under their  
127 hypothesis (such as 6–8–12), so failing to adequately challenge and disconfirm this  
128 hypothesis. On a constructivist perspective, learners can only base their exploration on  
129 testing hypotheses they have actually generated (or else behave randomly). To the extent  
130 that certain simpler hypotheses like “ascending numbers” were less likely to be generated  
131 on the basis of the provided example (cf. Oaksford & Chater, 1994; Tenenbaum, 1999), it is  
132 not surprising that participants failed to actively exclude these possibilities with their tests.

133 In the computational cognitive science literature, recent symbolic search ideas  
134 manifest under the label of “learning as program induction”. Such models have begun to be  
135 applied to synthesizing humanlike problem solving and planning and tool use (Allen,  
136 Smith, & Tenenbaum, 2020; Ellis et al., 2020; Lai & Gershman, 2021; Lake, Ullman,  
137 Tenenbaum, & Gershman, 2017; Ruis, Andreas, Baroni, Bouchacourt, & Lake, 2020; Rule,  
138 Schulz, Piantadosi, & Tenenbaum, 2018). We will draw on these in examining children and  
139 adults hypothesis generation.

<sup>140</sup> **Constructivism in Development**

<sup>141</sup> The “child as scientist” (Carey, 1985; Gopnik, 1996)—or recently, “child as hacker”  
<sup>142</sup> (Rule, Tenenbaum, & Piantadosi, 2020) — perspective casts children’s cognition as driven  
<sup>143</sup> by broadly the same inductive processes as adults’ but at an earlier stage in a journey of  
<sup>144</sup> construction and discovery.

<sup>145</sup> While children have been shown to be capable active learners (McCormack,  
<sup>146</sup> Bramley, Frosch, Patrick, & Lagnado, 2016; Meng, Bramley, & Xu, 2018; Sobel & Kushnir,  
<sup>147</sup> 2006) there is also evidence that children’s ability to learn effectively from active learning  
<sup>148</sup> data is more fragile than adults’. For example, children’s play can look repetitive and  
<sup>149</sup> inefficient when held to information theoretic norms (Lapidow & Walker, 2020; McCormack  
<sup>150</sup> et al., 2016; Meng et al., 2018; Sim & Xu, 2017). Sobel and Kushnir (2006) also found  
<sup>151</sup> children were much less accurate at causal structure identification in “yoked”  
<sup>152</sup> conditions—where they had to use evidence generated by someone else to learn—while  
<sup>153</sup> adults are less effected, sometimes able to learn about as well from others’ data as their  
<sup>154</sup> own (Lagnado & Sloman, 2006). This performance gap has been argued to stem from the  
<sup>155</sup> mismatch between whatever idiosyncratic hypotheses are under consideration by the  
<sup>156</sup> observer and those being tested by the active learner, making the yoked learner less able to  
<sup>157</sup> use the data to progress their theories (Fränken et al., 2022; Markant & Gureckis, 2014).  
<sup>158</sup> Relatedly, children have been argued to be more narrowly focused toward testing a single  
<sup>159</sup> hypothesis at a time (Bramley, Jones, Gureckis, & Ruggeri, 2022; Ruggeri & Lombrozo,  
<sup>160</sup> 2014; Ruggeri, Lombrozo, Griffiths, & Xu, 2016). This might reflect a less developed  
<sup>161</sup> working memory, restricting the number of hypotheses children can keep track of and  
<sup>162</sup> compare to evidence. An early emphasis on exploration has also been argued to be an  
<sup>163</sup> effective solution to a lifelong explore–exploit tradeoff, since earlier discoveries can be  
<sup>164</sup> exploited for longer (Gopnik, 2020). Program induction also provides a potential  
<sup>165</sup> explanation for transitions between developmental “stages”, characterized by occasional  
<sup>166</sup> leaps forward in insight. For instance, Piantadosi, Tenenbaum, and Goodman (2012)  
<sup>167</sup> demonstrate how a program induction model can reproduce a characteristic developmental  
<sup>168</sup> transition from grasping a few small numbers to discovering a recursive concept of real  
<sup>169</sup> numbers. We note that an important part of constructivism is the idea that we *cache* the  
<sup>170</sup> useful concepts we invent (cf. Zhao, Bramley, & Lucas, 2022), meaning our conceptual  
<sup>171</sup> library grows as we do, becoming richer and more powerful for solving the tasks we  
<sup>172</sup> repeatedly face. We do not attempt to model this important aspect of constructivism in  
<sup>173</sup> this paper but return to it in the General Discussion.

<sup>174</sup> Differences between childlike and adultlike inductive inference might also be  
<sup>175</sup> captured by parameterizable differences in search, potentially reflecting principles of

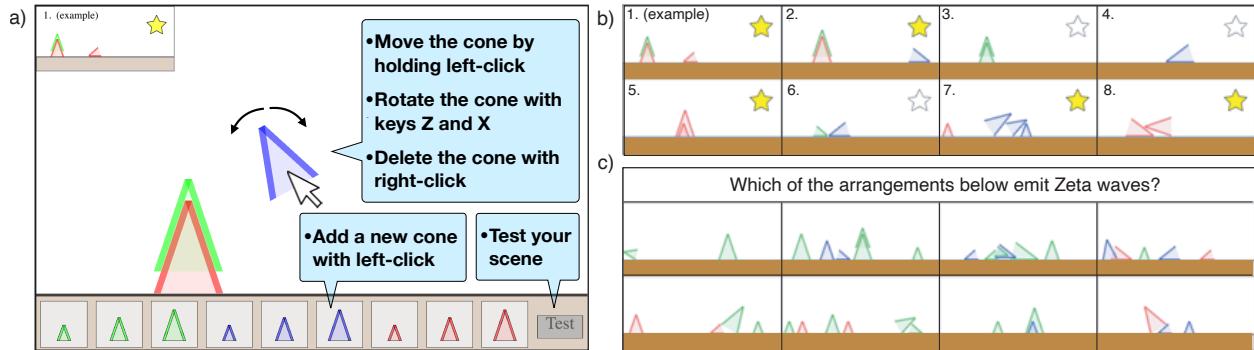
176 stochastic optimization (Lucas, Bridgers, Griffiths, & Gopnik, 2014). For instance, young  
177 children have been found to be quick to make broad abductive generalizations from a small  
178 number of examples—e.g. readily imputing novel physical laws to explain surprising  
179 evidence (L. E. Schulz, Goodman, Tenenbaum, & Jenkins, 2008). Building on this finding,  
180 children’s hypothesis generation and search has been framed as rationally “higher  
181 temperature” than adults’—producing more diversity of ideas at the cost of being noisier  
182 (Lucas et al., 2014). This is algorithmically sensible as optimization over high dimensional  
183 spaces is known to be more effective when proposals are initially large leaps and decrease  
184 over time, as in *simulated annealing* (Van Laarhoven & Aarts, 1987). However, a high  
185 diversity of guesses might also reflect that children have a rationally flatter latent prior  
186 than adults, inherently entertaining a wider range of hypotheses at the cost of entertaining  
187 high probability ones less frequently. A third possibility is that children’s hypothesis  
188 generation might be driven more by *bottom-up* processing than adults’. With less  
189 established expectations, or less powerful primitive concepts to work with, children’s  
190 hypotheses might more directly *describe* encountered patterns, while adults might rely  
191 more on their existing knowledge hierarchy to constrain hypothesis generation in a  
192 *top-down* way (Clark, 2012). We will contrast children’s and adults’ hypothesis generation  
193 and active learning in a rich task setting that allows us to closely investigate these ideas.

## 194 Task

195 In order to study inductive learning, we use a rich open-ended task that extends on  
196 Wason (1960) and the logical rule-induction tasks studied by Nosofsky et al. (1994), Lewis  
197 et al. (2014), Goodman et al. (2008), and Piantadosi et al. (2016). Akin to the  
198 blicket-detector paradigm in developmental causal cognition (Gopnik et al., 2004; Lucas et  
199 al., 2014), our task has a causal framing, probing inductive inferences about what  
200 conditions make an effect occur in a minimally contextualized domain. However, departing  
201 from Blicket detector tasks, we include a large and physically rich set of features that  
202 learners can draw on in their inferences allowing test scenes to vary in the number, nature  
203 and arrangement of objects. Our task is inspired by a tabletop game of scientific induction  
204 called “Zendo” (Heath, 2004) and builds on a pilot task examined in (Bramley, Rothe,  
205 Tenenbaum, Xu, & Gureckis, 2018). In it, learners both observe and create *scenes*, which  
206 are arrangements of 2D triangular objects called *cones* (Figure 1) and test them to see if  
207 they produce a causal effect (which arrangements of blocks “make stars come out” in our  
208 minimal framing). The goal is to both predict which of a set of new scenes will produce the  
209 effect and describe the hidden rule that determines the general set of circumstances  
210 produce the effect (try it [here](#)). Scenes could contain between 1 and 9 cones. Each cone has

211 two immutable properties: size $\in\{\text{small, medium, large}\}$  and color $\in\{\text{red, green, blue}\}$  and  
 212 continuous scene-specific  $x\in(0,8)$ ,  $y\in(0,6)$  positions and orientations $\in(0,2\pi)$ . In addition to  
 213 cones' individual properties, scenes also admit many relational properties arising from the  
 214 relative features and arrangement of different cones. For instance, subsets of cones might  
 215 share a feature value (i.e., be the same color, or have the same orientation) or be ordered  
 216 on another (i.e., be larger than, or above) and pairs of cones might have relational  
 217 properties like pointing at one another or touching. This results in an extremely rich  
 218 implicit space of potential concepts.

219 We note that, by design, the dimensionality of this task makes it extremely difficult.  
 220 As with Wason's 2-4-6 example, and genuine questions of scientific induction, the hard part  
 221 of this task is not evaluating whether a candidate hypothesis can explain the data but  
 222 rather generating the right hypothesis in the first place. As with the 2-4-6 task, there are  
 223 always infinite data-consistent possibilities and while the bulk of these may be outlandishly  
 224 complex, many others may still be simpler or more salient than the ground truth. Without  
 225 carefully gathered evidence with broad coverage of the space of possible scenes, a learner  
 226 will frequently be unable to rule out simpler possibilities that more parsimoniously capture  
 227 the data than the ground truth, essentially being left with evidence that would not lead  
 228 even an unbounded Bayesian agent to the correct answer.<sup>1</sup>



**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.

229 We use mixed-methods (Johnson, Onwuegbuzie, & Turner, 2007), analyzing both

<sup>1</sup> In tabletop game form, Zendo typically takes dozens of rounds of tests and incorrect guesses by multiple guessers, as well as leading examples and clues from the rule-setter for even simple hidden rules to be identified. An online community on Reddit play a binary sequence version of Zendo, often taking hundreds of guesses before the answer is found if it is at all (for example [here](#)).

230 qualitative data in the form of freely generated guesses about the symbolic rules and  
 231 quantitative data in the form of forced choice generalizations. Concretely, we adopt an  
 232 expressive concept grammar inspired by constructivist ideas in developmental psychology  
 233 and formalized using program induction ideas from machine learning. We assume the  
 234 latent space of possible concepts in our task are those expressible in first order logic  
 235 combined with lambda abstraction (Church, 1932) and full knowledge of the potentially  
 236 relevant features of the scene (see Appendix Table A-1 for the grammatical primitives we  
 237 assume). Table 1 shows the five ground truth rules we used in our experiment expressed in  
 238 natural language and in lambda calculus along with the initial rule-following example scene  
 239 we provided to participants.

240 Given the inherent difficulty of this type of task we expect absolute accuracy to be  
 241 fairly low for both children and adults (and for our models). However, we expect that  
 242 many participants will be able to make guesses that are consistent with most of the  
 243 evidence they have. Since we might expect evaluation of evidence–hypothesis consistency  
 244 to be more error-prone in children, we expect adults’ guesses to be more strictly consistent  
 245 with their evidence. Finally, there is the question of relative dominance of bottom-up and  
 246 top-down processing in children’s and adults’ guesses. To explore this, we consider two  
 247 models that differ in this dimension.

#### 248 Context-free hypothesis generation

249 In examining children’s and adults’ inferences, we start by laying out a “top-down  
 250 first” approach to hypothesis generation, utilizing a probabilistic context-free grammar  
 251 (PCFG) to define and draw from a latent prior over concepts expressible in first order  
 252 logic. A PCFG is a collection of “construction rules” that, when run repeatedly,  
 253 stochastically create expressions in an underlying grammar (Ginsburg, 1966). A PCFG can  
 254 be used to generate a prior sample of hypotheses that can then be weighted by their  
 255 likelihoods of producing observations—here, their ability to reproduce the labels of the  
 256 scenes that the participant has tested. The hypotheses make predictions about new scenes  
 257 which can be weighted by their posterior probability and marginalized over to make  
 258 generalizations. Because parts of this production process and underlying grammar involve  
 259 branching—e.g., “and” and “or”—sampled hypotheses can be arbitrarily long and complex,  
 260 involving multiple Boolean functions and complex relationships between an unlimited  
 261 number of bound variables. In this way, an infinite latent space (in our case first order logic  
 262 + lambda abstraction) is covered in the limit of infinite PCFG sampling (see Figure 2a).  
 263 Thus, one way to think of the PCFG is as a *computational level* characterization of the  
 264 problem of inductive inference. However, we will argue that the generative mechanism at

265 the heart of of the PCFG framework also elucidates important mechanistic considerations  
 266 and provides the representational framework needed to ground algorithmic approximations  
 267 that depart from this ideal and reflect core constructivist ideas.

268 At the computational level, different PCFGs, containing different primitives and  
 269 expansions, can be compared against human behavior. And the probabilities for the  
 270 productions in a PCFG can be fit to maximize correspondence with human judgments. In  
 271 this way, recent work has attempted to infer the “logical primitives of thought” (Goodman  
 272 et al., 2008; Piantadosi et al., 2016). Here we consider a single expressive PCFG  
 273 architecture and examine its behavior under limited sampling. We examine its behavior  
 274 with uniform production weights but also with weights engineered to produce the  
 275 characteristics of “childlike” and “adultlike” symbolic guesses in our task. Crucially, under  
 276 all these weighting schemes, our PCFG embodies the principle of parsimony: Simpler  
 277 concepts—composed of fewer grammatical parts (Feldman, 2000)—have a higher  
 278 probability of being produced and so are favored over more complex ones equally able to  
 279 explain the data.

280 While naively, we might expect children to entertain simpler concepts than adults,  
 281 this induction framework tends to predict the reverse. If we assume we start life at our  
 282 most flexible, or “programable” (Turing, 2009), this would be like being born with concept  
 283 building mechanism that is initially “untuned”, growing its concepts essentially through  
 284 blind mutation (Campbell, 1960) where each forking path on the road to a complete  
 285 concept starts out equiprobable. However as a learner gathers a lifetime of experience, we  
 286 would expect these construction weights to become tuned so as to favor certain elements or  
 287 features that have proven useful in the past. A uniform-weighted PCFG hypothesis  
 288 generator will thus tend to produce greater diversity than a more fine-tuned one. As such,  
 289 it embodies the idea that more elaborately or implausibly structured, or “weird”, concepts  
 290 will come to the minds of children than adults.

291 What PCFG approaches have in common is a generative mechanism for sampling  
 292 from an infinite latent prior, here over possible logical concepts. However, sampled  
 293 “guesses” must also be tested against data. Unfortunately, in our task—and perhaps even  
 294 more so outside of it—the vast majority a priori generated concepts are likely to be  
 295 inconsistent with whatever evidence a learner has already encountered.<sup>2</sup> For this reason,  
 296 the procedure is astronomically inefficient, requiring very large numbers of samples in order

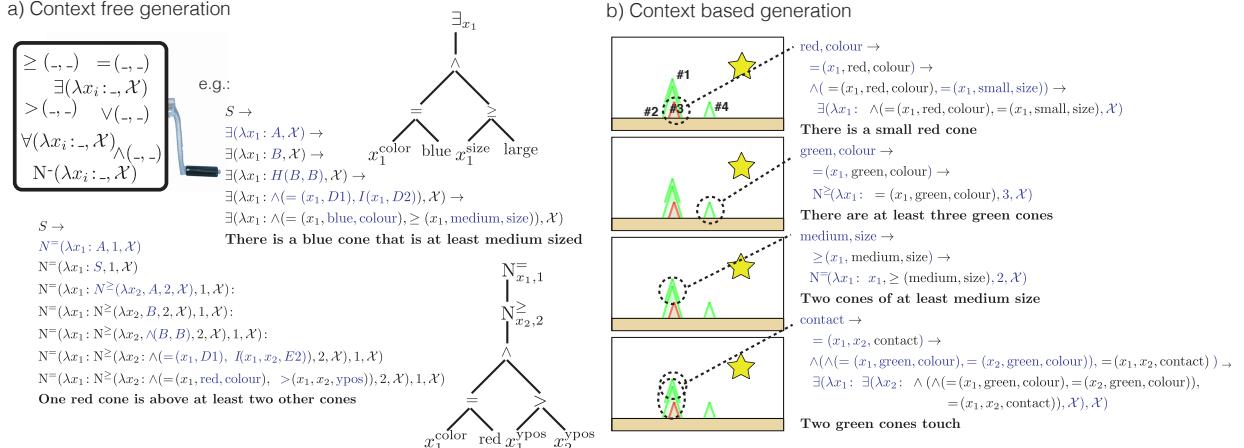
---

<sup>2</sup> In our task, mere simply tautological (i.e., “All cones are red or not red”), contradictory (i.e., “There is a cone that is red and not red”), physically impossible (“Two (different) objects have the same position”) Indeed, around 20% of the hypotheses generated by our PCFGs are tautologies, and 15% are contradictions. Many others combine a meaningful hypothesis with a tautological corollary (i.e., “There is a large red object that is larger than all medium sized objects”).

297 to reliably generate non-trivial rules. One can also use a PCFG to adapt existing  
 298 hypotheses, for instance using a Markov Chain Monte Carlo scheme in which parts of a  
 299 hypothesis are regrown and accepted according to their fit to evidence (cf. Fränken et al.,  
 300 2022; Goodman et al., 2008). While we think this approach is promising we do not model  
 301 this here, and simply return to it in the general discussion. However, we do additionally  
 302 consider an alternative to the PCFG, that provides a more sample efficient and, on the face  
 303 of it, more cognitively plausible mechanism for initializing new hypotheses.

### 304 Context-based hypothesis generation

305 Instance Driven Generation (IDG) (Bramley et al., 2018) is a recent proposal  
 306 related to the PCFG framework but with a key difference. Rather than generating initial  
 307 hypotheses prior to, or blind to the current evidence, the IDG generates ideas *inspired* by  
 308 encountered patterns (cf. Michalski, 1969), thus incorporating bottom-up reactivity to  
 309 evidence into its conceptualization process. Each IDG hypothesis starts with an  
 310 observation of features of one or several objects in a scene and uses these to back out a true  
 311 logical statement about the scene in a stochastic but truth-preserving way. If the scene is  
 312 rule following, this statement constitutes a positive hypothesis about the hidden rule.  
 313 Otherwise, it constitutes a negative hypothesis, i.e. about what must *not* be present. Thus,  
 314 an IDG does not begin each learning problem with a prior over all possible concepts, but  
 315 rather draws its initial ideas from a restricted space consistent with the extant patterns in  
 316 a focal observation. Figure 2b illustrates this approach. While a regular PCFG effectively  
 317 starts at the top level (i.e. outermost nesting) of a compound concept and works downward  
 318 and inward, the IDG starts from the central content (drawn from its observation) and  
 319 works upward and outward to a quantified statement, ensuring at each step that the  
 320 statement is true of the scene. The result is a mechanism that uses its concept grammar to  
 321 describe features and patterns in evidence. This means that the IDG does not entertain  
 322 hypotheses that are possible but never exemplified by a scene. For example, “at most five  
 323 reds” would only be generated if a learner actually saw a rule-following scene containing  
 324 five reds. A key prediction of the IDG is an interaction between the scenes generated by  
 325 the participant and the hypotheses these subsequently inspire, with simpler scenes,  
 326 embodying fewer extraneous or coincidental patterns being more likely to inspire the  
 327 learner to generate the true concepts.

**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 A.

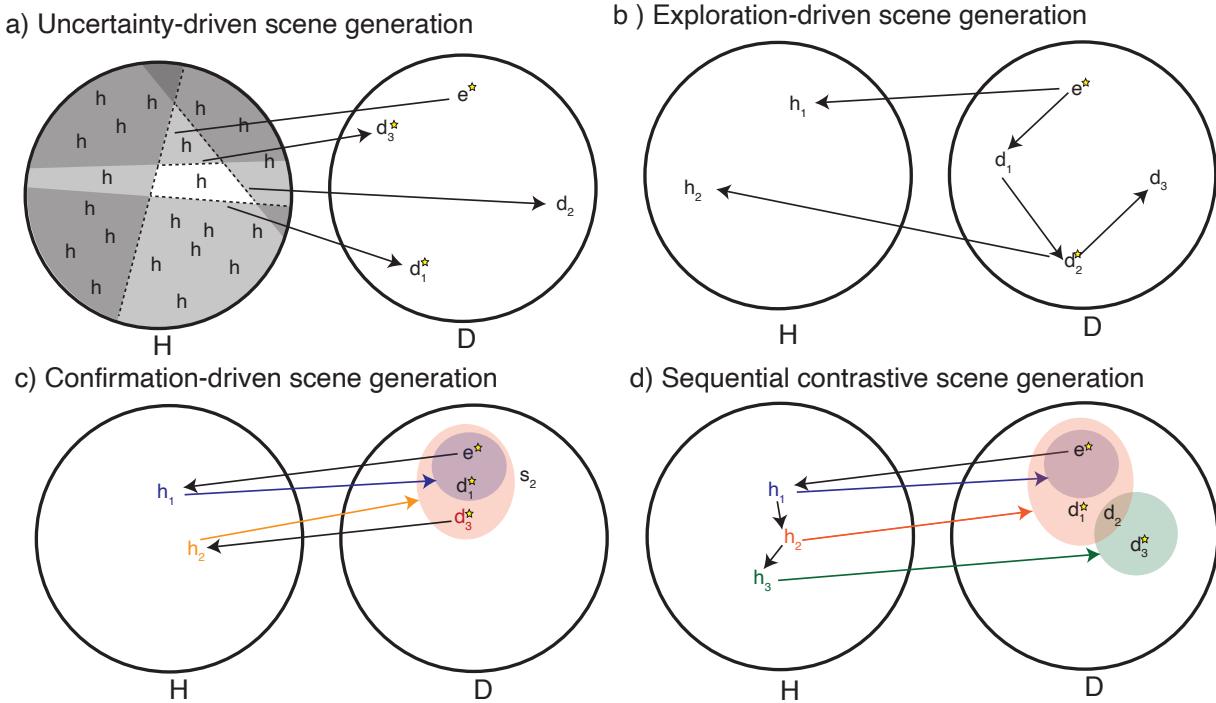
### 328 Hypothesis-driven scene generation

#### 329 Uncertainty-driven learning

330 Normatively, test scenes should serve to minimize expected uncertainty across the  
 331 full hypothesis space. A direct way to approximate this here is to start with a prior sample  
 332 of hypotheses (e.g. drawn context-free) and progressively create scenes that serve to  
 333 minimize expected uncertainty over this sample by forking their predictions (Bramley et  
 334 al., 2022; Nelson, Divjak, Gudmundsdottir, Martignon, & Meder, 2014). We visualize this  
 335 in Figure 3a, imagining three labelled scenes  $d_1 \dots d_3$  that progressively divide a prior  
 336 sample of hypotheses ( $hs$ ) until a most-likely candidate emerges. The constructivist setting  
 337 presents a challenge for this norm since the hypothesis space is latent and is initially  
 338 unexplored.

#### 339 Exploration-driven learning

340 An alternative hypothesis-free approach might be to explore the data space directly,  
 341 for instance generating scenes that vary in the number and nature of objects they contain  
 342 in the hope of naturally uncovering concept boundaries and inspiring hypothesis  
 343 generation. We sketch this in Figure 3b. Efficient uncertainty-driven and  
 344 exploration-driven learning both predict generation of scenes that differ substantially from  
 345 one another, ideally being anti-correlated so as to cover the space efficiently (Osborne et  
 346 al., 2012). However this does not seem well matched to constructivism, where we rather  
 347 think of the learner as entertaining a small but not completely empty set of possibilities

**Figure 3**

Active learning strategies:  $H$  = latent hypothesis space  $D$  = data space. Arrows indicate direction of inferences. a) Uncertainty-driven tests over prior sample  $h \in H$ . Dotted lines separate hypotheses by outcomes they predict for initial example  $e$  and self-generated scenes  $d_1 \dots d_3$ . Shading indicates which  $h$ s mis-predict each outcome. b) Exploration-driven testing. Scenes selected to explore  $D$  without regard to  $H$ . c) Confirmatory testing: Example  $e$  inspires hypothesis  $h_1$ . Scenes then test its generalization predictions. Colored circles visualize space of scenes for which each hypothesis predicts outcome will be produced.  $d_1$  and  $d_2$  are correctly predicted as rule following.  $d_3$  is mispredicted by  $h_1$  in producing the outcome, leading to a new  $h_2$ . d) Sequential contrastive testing:  $e$  inspires  $h_1$  and  $h_1$  inspires  $h_2$ ,  $d_1$  contrasts these leading to rejection of  $h_1$ .  $h_2$  then inspires  $h_3$  and  $d_2$  contrasts these, etc.

348 and hence unable to capitalize on such diverse evidence.

349 A constructivist way to think of active learning is as acting in ways that challenge  
 350 one's current hypotheses and so facilitate their refinement or the construction of better  
 351 alternatives. We sketch two such approaches: Confirmatory testing and Sequential  
 352 Contrastive testing.

### 353 **Confirmatory testing**

354 With a candidate hypothesis in mind, a learner can seek to challenge it through its  
 355 generalizations (Nickerson, 1998; Popper, 1959). For example, after encountering the scene  
 356 in row 1 of Table 1, a learner might generate the initial hypothesis that “there must be a

357 small red” (since this describes one of the objects). To confirm this, they might try a  
 358 positive generalization test, i.e. keep the small red but remove or randomize the other  
 359 objects and predict the effect will still occur (e.g.  $d_1$  in Figure 3c). Alternatively they  
 360 might use it to predict a way to minimally alter  $d_1$  so it no longer produces the effect,  
 361 removing the small red and keeping the rest (e.g.  $d_2$ ). So long as the learner gets the  
 362 outcome they anticipate, they can stick with their hypothesis. When they don’t they can  
 363 either abandon or adapt it. For instance,  $d_3$  in Figure 3c proves inconsistent with  $h_1$ ,  
 364 requiring a new hypothesis be generated that can explain why  $d_1$  and  $d_3$  produce the effect  
 365 but not  $d_2$ . A limitation of a one-hypothesis-at-a-time approach is that it is unclear how  
 366 distinctive the hypothesis’s generalization predictions are.<sup>3</sup> For example, since the ground  
 367 truth in this example is just “there is a red”, producing new scenes containing small reds  
 368 will fail to reveal that the redness but not the smallness is causative of the label. Another  
 369 limitation is that it is unclear what to do when one’s hypothesis is ruled out, especially if  
 370 the scene if the test that differs dramatically from the ones with which it is consistent. For  
 371 this reason, the education literature has long emphasized the utility of a “*control of*  
 372 *variables*” strategy (Chen & Klahr, 1999; Klahr, Fay, & Dunbar, 1993; Klahr, Zimmerman,  
 373 & Jirout, 2011). This amounts to manipulating exactly one design variable per test, such  
 374 that any difference in the outcome is straightforwardly attributable to the change in the  
 375 input providing a route to adapting one’s hypothesis when it fails.

376 ***Sequential contrastive testing***

377 A related scheme that might allow a constructivist learner to escape some  
 378 pathologies of confirmatory testing is the *iterative counterfactual strategy* described in  
 379 Oaksford and Chater (1994). That is, learners might first generate an *alternative*  
 380 *hypothesis*  $h_2$  by inverting some feature of their initial hypothesis and then focus their next  
 381 test on separating  $h_1$  from  $h_2$  (e.g., Figure 3d).<sup>4</sup> For example, starting with  $h_1$ : “there is a  
 382 small red”, one local alternative would be to drop the the mention of size, leading to  $h_2$ :  
 383 “There is a red”. Now the learner has a pair of hypotheses and a recipe distinguishing  
 384 between them: Testing a scene containing a red object that is not small (e.g.  $d_1$ ). This  
 385 could again be easily achieved by adapting the original scene, so the small red is a different

---

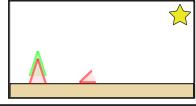
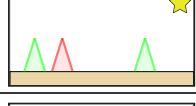
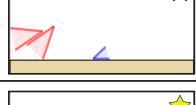
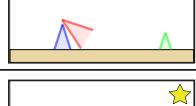
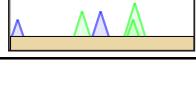
<sup>3</sup> A general finding is that positive confirmatory tests are valuable to the extent that the outcome of interest is rare, e.g. if most scenes are not rule following. This is not generally the case in this task.

<sup>4</sup> In Oaksford and Chater’s (1994) formulation, the complementary hypothesis is then inconsistent with the scene that inspired the original hypothesis, such as going from “increasing by two” (inspired by seeing 2-4-6) to “decreasing by two” such that its falsification may be mistaken for confirmation of the original hypothesis. Here there are many ways to flip the content of a hypothesis both with or without rendering it inconsistent with a scene that inspired it.

size (Chen & Klahr, 1999; Klahr et al., 1993, 2011). If  $d_2$  produces the effect,  $h_1$  can be supplanted with  $h_2$ . Otherwise  $h_2$  can be rejected and a new  $h_3$  can be generated. Either way, this approach facilitates constructivism by providing a direction of travel however a test comes out, so allowing a constructivist learner to explore both the data and hypothesis spaces in parallel (Klahr & Dunbar, 1988).

As illustrated in Figure 3, what constructivism-compatible hypothesis-driven approaches have in common is a prediction of anchoring in data space: Each new scene shares features with the scene that inspired the earlier hypotheses that inspired it. This contrasts with the pattern we would expect if participants followed a normative uncertainty-driven approach or model-free exploration-driven approach since both naturally predict each scene should be as different as possible to earlier ones. While we do not collect the trial-by-trial guesses we would need to distinguish between all the accounts we mention, we will look for an empirical signature of constructivist active learning, in the form of anchored, incremental and systematic testing patterns and assess whether these differ between children and adults.

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

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

## 401 Overview

In summary, the main goal of this paper is a close investigation of developmental differences in active open-ended hypothesis generation examined through the lens of a constructivism-inspired rational-process framework that puts stochastic generation and

incremental search at the center of the individuals' learning. To foreshadow, we find that children make more complex guesses about the hidden rule that are only a marginally worse fit to the evidence than adults' guesses. Children also create more complex learning data than adults but do so less systematically. We then show that both children's and adults' guesses reflect an evidence-inspired process of compositional concept formation over a top-down-first PCFG norm, capturing that their guesses are inspired by discovery of patterns in their learning data. We show these behavioural patterns are a natural result of children having a less fine-tuned concept generation mechanism. Crucially, we also show that both children's and adults' symbolic guesses causally drive their generalizations, as opposed to these being driven by surface feature resemblance as emphasized in statistical views of concepts (cf. Medin & Schaffer, 1978; Posner & Keele, 1968). Finally, we show that both children's and adults' create scenes by adapting earlier scenes, which we argue is consistent with confirmatory or iterative counterfactual testing rather than uncertainty- or exploration-driven testing.

**419 Experiment**

**420 Methods**

**421 Participants**

We recruited 54 children in the lab (23 female, aged  $8.97 \pm 1.11$ ) and 50 adults online (22 female, aged  $38.6 \pm 10.2$ ). Forty children completed all five trials and the remaining 14 completed  $2.71 \pm 1.07$  trials before indicating that they had had enough. For these children we simply include the trials that they completed. We collected participants until we reached our intended sample size of 50 per agegroup after exclusions. We chose this sample size simply to exceed our 2018 ( $N=30$ ) pilot with adults.<sup>5</sup> Ten additional adult participants completed the task but were excluded before analysis for providing nonsensical or copy-pasted text responses. Adult participants were paid \$1.50 and a performance related bonus of up to \$4 ( $\$1.96 \pm 0.75$ ). Children's sessions lasted between 30 minutes and an hour. For adults, the task took  $27.49 \pm 12.09$  minutes of which  $9.8 \pm 7.9$  was spent on instructions. The children's and adults' versions of the task are available to try here [https://github.com/bramleyccslab/computational\\_constructivism](https://github.com/bramleyccslab/computational_constructivism).

---

<sup>5</sup> While we note that 104 is not a large sample by modern standards, our focus is on modeling inferences at the individual level. Each participant produces an exceptionally rich dataset and our analyses have unusually large storage and compute requirements making a larger sample infeasible to analyze.

434 ***Design***

435 All participants faced the same five learning problems in an independently  
436 randomized order (see Table 1). For each learning problem participants were given an  
437 initial positive example, as shown in the table, and then performed self tests of their own  
438 before making generalizations and free guesses as to the hidden rule.

439 ***Materials and Procedure***

440 **Child sample.**

441 **Instructions.** Participants sat in front of a laptop with a mouse attached, with  
442 the experimenter sitting next to them and interacted with the task through the browser.

443 The experimenter read out the instructions for the participant. These explained  
444 how the game worked and showed the participant five examples of possible rules the blocks  
445 could have (relating to color, size, proximity, angle, or relation). The instructions also  
446 included videos showing the participant how to manipulate the blocks using the mouse and  
447 keyboard. After the instructions, the participant was given a comprehension check of five  
448 true or false questions. If they did not get them all right on their first try, the experimenter  
449 read through the instructions again and asked them again. All participants passed the  
450 comprehension check the second time.

451 **Learning Phase.** The participant was then introduced to an initial example of a  
452 block type (“Here are some blocks called [name]s. We’re going to click test to see if stars  
453 will come out of the [name]s.”). The initial example of each block type (i.e., each rule) was  
454 constant across participants. Since every initial example of a block type was a positive  
455 example, a star animation played when the “Test” button was clicked. The participant was  
456 encouraged to use either the trackpad or the mouse to click the “Test” button, whichever  
457 was comfortable for them.

458 After the initial positive example, the participant was shown a blank scene with  
459 blocks available to add to it, and was asked to test the blocks seven more times  
460 (Figure 1a). The scene creation interface was subject to simulated gravity, meaning there  
461 were physical constraints on how the objects can be arranged. The experimenter told them  
462 they could now play with the blocks like they saw in the instructional video. The  
463 experimenter also reminded the participant of how to add, remove, move, and rotate blocks  
464 on the screen using the mouse and keyboard. Participants were encouraged to ask for help  
465 with moving the blocks if needed. If they seemed to be having trouble, the experimenter  
466 would ask if they needed help with setting up the blocks. The participants were told that  
467 when they had finished moving the blocks around, they should press the “Test” button to  
468 see if stars came out of them. For positive tests, the experimenter would neutrally say:

469 “Stars did come out of the [name]s that time” and for negative tests: “Stars did not come  
470 out of the [name]s that time.”

471 **Question Phase.** After testing the blocks a total of eight times (Figure 1b),  
472 participants were shown a selection of eight more pre-determined scenes containing blocks  
473 (Figure 1c). The experimenter asked them to click on which pictures they thought the  
474 stars would come out of, reminding them that they could pick as many as they wanted, but  
475 they had to pick at least one. Unknown to participants, half of these scenes were always  
476 rule following but their positions on screen were independently counterbalanced. The test  
477 scenes and their labels remained visible on the screen throughout the Learning and  
478 Question phases.

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

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

487 **Adult sample.** We recruited our adult sample from Amazon Mechanical Turk  
488 and adults completed the task on their own computers. They completed the same  
489 instructions as the children with an additional section about bonuses and had to  
490 successfully answer comprehension questions, including an additional two about the  
491 bonuses, before starting the main task. Specifically, adults were bonused 5 cents for each  
492 correct generalization (up to a possible 40 cents for each of the five trials) and an  
493 additional 40 cents for a correct guess as to the hidden rule, again for each of the five trials.  
494 Aside from having no experimenter in the room, and filling out the text fields themselves,  
495 the procedure was identical to the children’s task. Full materials including experiment  
496 demos, data and code are available at the [Online Repository](#).

## 497 Results

498 We first look at the qualitative characteristics of children’s and adults’ explicit rule  
499 guesses then assess relative accuracy of participants’ rules and generalizations about new  
500 scenes before comparing the features of the scenes produced by adults and children. We  
501 will then turn to a series of model-based analyses that attempt to reproduce participants  
502 distributions of free guesses, generalizations and scenes within the constructivist framework.

503 ***Guess complexity and constituents***

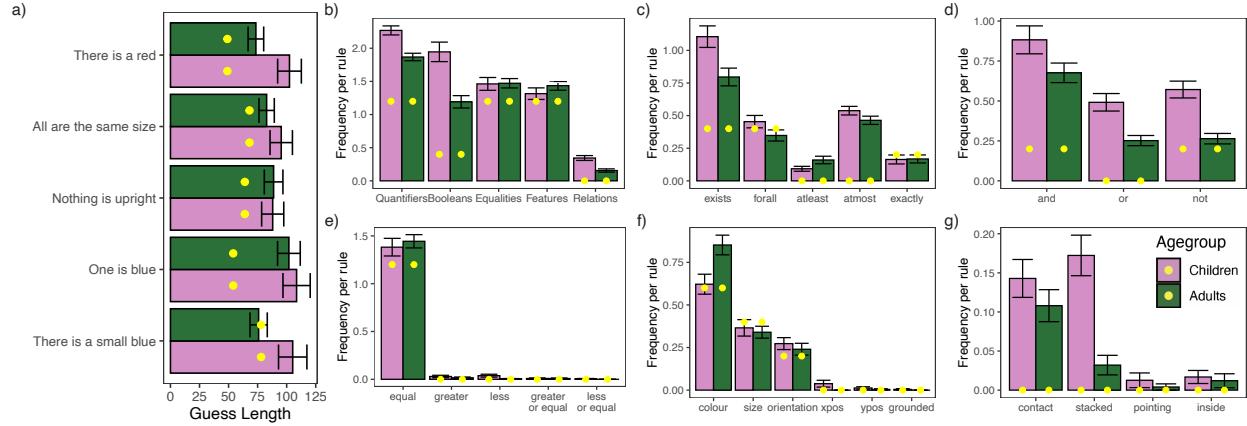
504 We had human coders translate participants' free text guesses about the hidden rule  
 505 wherever possible into an equivalent logical expression using the grammatical elements  
 506 available to our learning models. We were able to do this for 86% (n=205) of children's  
 507 trials and 88% (n=219) of adults' trials. For example, if the participant wrote "*There must*  
 508 *be one big red block*" this was converted into

509  $N^=(\lambda x_1 : \wedge(=(x_1, \text{large}, \text{size}), =(x_1, \text{red}, \text{color})), 1, \mathcal{X})$ . This logical version can be  
 510 automatically evaluated on the scenes and can be read literally as asserting "*There exists*  
 511 *exactly one  $x_1$  in the set of objects  $\mathcal{X}$  such that  $x_1$  has the size 'large' and the color 'red'*".  
 512 We had a primary coder, blind to the experimental hypotheses code all responses, and a  
 513 second coder blind spot check 15% of these (64). The two coders agreed in 95% of cases.  
 514 We provide further details about the coding in Appendix B and full coding resources and  
 515 full coding data in the [Online Repository](#).

516 To explore structural differences in children's versus adults' hypotheses, we first  
 517 break down these encoded rule guesses into their logical parts. This primarily reveals that  
 518 children's encoded rules were substantially *more complex* than those generated by adults  
 519 and that both were substantially more complex than the ground truth rules. Children's  
 520 and adults' rules also differed in terms of the prevalence of particular elements and features  
 521 (see Figure 4). As an example, one child's rule for problem 1 was "*You must have two reds*  
 522 *and one blue*" which was translated to

523  $N^=(\lambda x_1 : N^=(\lambda x_2 : (\wedge(=(x_1, \text{red}, \text{color}), =(x_2, \text{blue}, \text{color})), 1, \mathcal{X}), 2, \mathcal{X})$ , requiring two  
 524 quantifiers ( $N^=$ ), one boolean ( $\wedge$ ), 2 equalities ( $=()$ ), and two references to the feature  
 525 color. The typical child-generated-rule used 2.25 quantifiers (4c), 2.06 booleans (4d), 1.55  
 526 equalities and inequalities (4e), referred to 1.39 different primary features (color, size,  
 527 orientation, x- or y-position, groundedness, 4f) and 0.37 relational features (contact,  
 528 stackedness, pointing, or insideness, 4g). In contrast, the average adult generated rule  
 529 required just 1.84 quantifiers, 1.20 booleans, 1.47 equalities and inequalities, and referred  
 530 to 1.44 primary features but only 0.16 relational features. Children thus used significantly  
 531 more quantification (i.e. referred to more separate entities)  $t(102) = 3.98, p < .0001$ , more  
 532 booleans  $t(102) = 3.59, p < .0001$  and relational features  $t(102) = 3.12, p < .002$  than  
 533 adults, but the agegroups did not differ significantly in mentions of (in)equalities  
 534  $t(102) = -0.05, p = 0.96$  and references to the objects' basic features  
 535  $t(102) = -.91, p = .36$ . When children posited that an "at least", "at most" or "exactly" a  
 536 certain number of objects must have certain features, the number they chose was  
 537 substantially higher than that for adults (2.36 compared to 1.58,  $t(68) = 3.72, p = 0.0004$ ).  
 538 In terms of features, adults frequently gave rules relating to color (58% compared to 39% of

539 children's rules,  $t(102) = 2.27, p = 0.025$ , while children were more likely to refer to  
 540 positional properties (26% compared to 18% of adults' rules  $t(102) = 2.15, p = 0.034$ ).



**Figure 4**

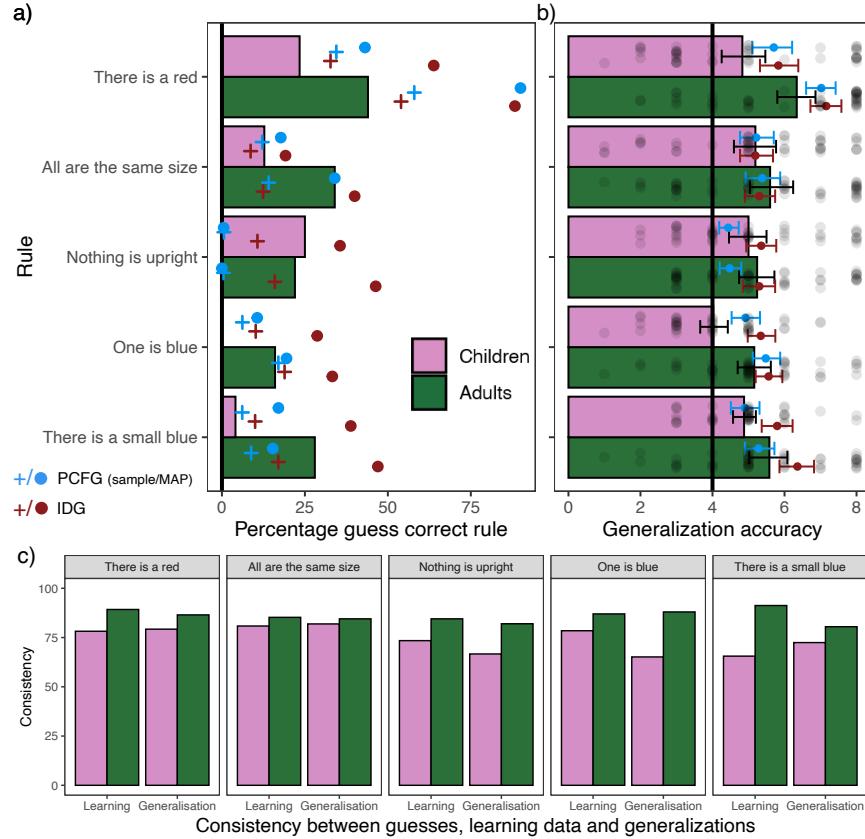
(a) Length of Children's and Adults' rule guesses. (b) Relative frequency of rule elements in logic coded versions of these rules, c–g with respect to quantifiers, booleans, (in)equalities, basic and relational features respectively. Error bars show normal 95% confidence intervals. Yellow points in a show ground truth frequency.

#### 541 Accuracy

542 Having observed systematic differences in the content of children's and adults'  
 543 hypotheses, we now ask if these manifest in children's and adults' inferential success; their  
 544 ability to identify the ground truth and make accurate generalizations.

545 **Guesses.** Both children and adults were occasionally able to guess exactly the  
 546 correct rules, doing so a respective 11% and 28% of trials. Adults produced the correct rule  
 547 more frequently than children  $t(102) = 4.0, p < .001$  and were more likely than children to  
 548 guess correctly (at a corrected significance level of 0.01) for the "All are the same size",  
 549 "One is blue" and "There is a small blue" rules (see Figure 5a). The plot reveals that no  
 550 child identified rule 4 exactly "One is blue" and only one identified rule 5 "There is a small  
 551 blue", while a slightly greater proportion of children than adults identified the positional  
 552 "Nothing is upright" rule. Note that chance level baseline for these free guesses is  
 553 essentially 0%. There are an unlimited number of wrong guesses and a small set of  
 554 semantically correct guesses. It is also the nature of this inductive problem that there are  
 555 an infinite number of wrong yet perfectly evidence-consistent rules for any evidence and  
 556 often there is a simpler evidence-consistent rule available than the ground truth.<sup>6</sup> Thus, it

<sup>6</sup> Although as more evidence arrives the ground truth is increasingly likely to be among "simplest" rules in a posterior sample.

**Figure 5**

a) Percentage children and adults guessing correct rule. b) Generalization accuracy. Bars show mean  $\pm$  bootstrapped 95% CIs. In a–b, Black vertical lines denote chance performance. Blue and red points show performance of simulated PCFG and IDG learners as described in Modeling section. Circles = guessing the MAP rule or MAP generalization (after marginalizing over posterior). “+” shows accuracy of a single posterior sample. Both models here use agegroup-consistent production weights, CIs show bootstrapped 95% confidence intervals. c) Consistency between subjects’ rule guess and their (self-generated) learning data, and generalizations.

557 is instructive to ask whether participants’ rules, where not exactly correct, are nevertheless  
 558 consistent with the evidence they gathered.

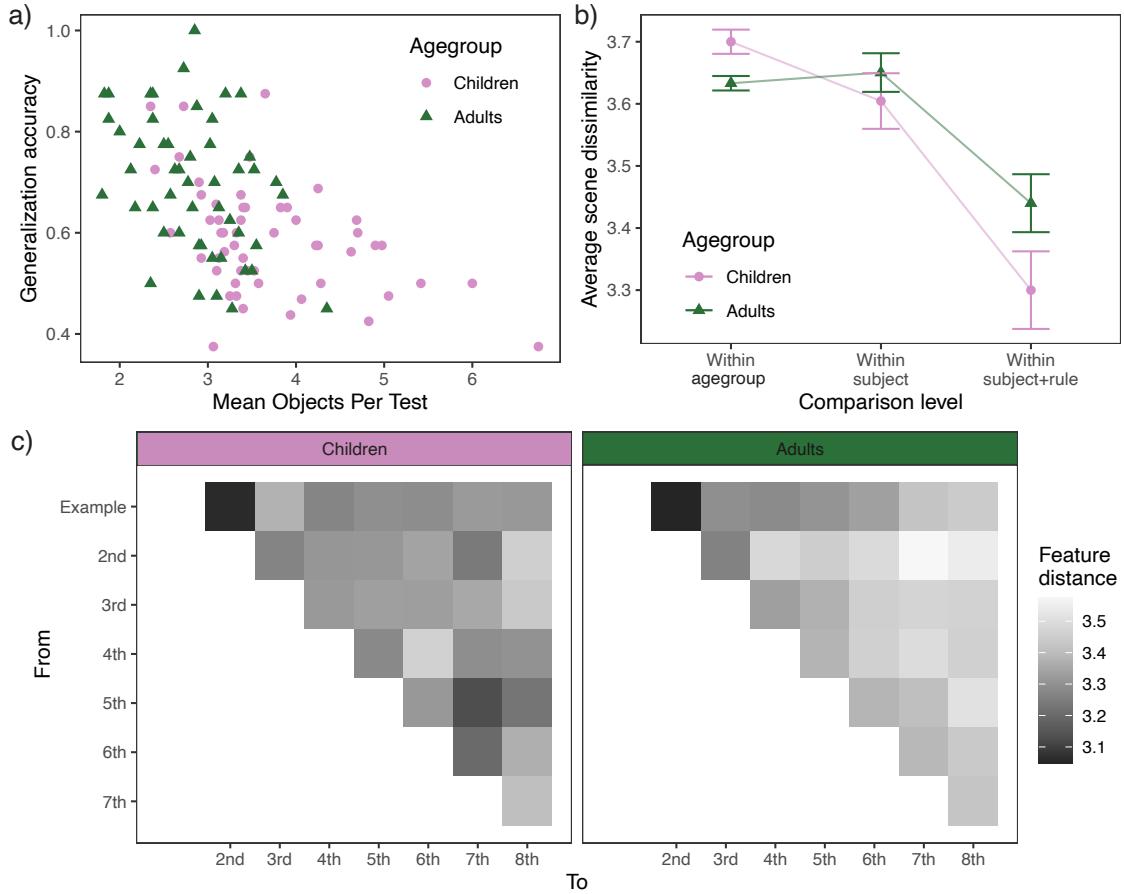
559 While, a completely random rule would only be consistent with all 8 scenes around  
 560  $0.5^8 \times 100 = 0.4\%$  of the time, children’s explicit rule guesses were perfectly consistent with  
 561 the labels of the 8 training scenes 30% of the time and Adult’s guesses were fully consistent  
 562 54% of the time. There was a moderate difference in average proportion of the learning  
 563 data explained by children’s compared to adults’ rules  $71\% \pm 27\%$  vs  $87\% \pm 17\%$   
 564  $t(98) = 5.6, p < .001$ . Similarly there was a difference the proportion of the participants’  
 565 generalizations that were consistent with their rule guess  $72\% \pm 21\%$  vs  $84\% \pm 16\%$ ,  
 566  $t(98) = 4.1, p < .001$  (see Figure 5c for a by-rule breakdown).

**Generalizations.** We now report participants performance in predicting which of 8 new scenes will produce stars (i.e. follow each hidden rule). Across the five tasks, both children and adults guessed more accurately than chance (50%): *children* mean $\pm SD$  59%  $\pm$  11%,  $t(53) = 5.9, p < .001$ ; *adults* 70%  $\pm$  14%,  $t(49) = 10.3, p < .001$ . Adults' generalizations were significantly more accurate than children's  $t(102) = 4.6, p < .001$  and children's accuracy improved significantly with age  $F(1, 52) = 6.2, \eta^2 = .11, p = 0.015$ . Indeed, adults' generalization accuracy was above a Bonferroni-corrected chance level of  $p \leq 0.01$  for all five rules and children were similarly above chance except for rules 1. "There is a red" ( $t(46) = 2.5, p = .015$ ) and 4. "One is blue" ( $t(46) = .1, p = .915$ ; see Figure 5b).

### 577 ***Scene generation***

As well as generating more complex rules, children tended to create more complex test scenes than adults. The average child-generated scene contained 3.7 $\pm$ 0.88 objects (close to the average in the example scenes) compared to 2.8 $\pm$ 0.57 objects for adults ( $t(102) = 5.8, p < .001$ ). The complexity of a learner's test scenes was inversely related to their performance overall ( $F(1, 102) = 39.0, \beta = -0.08, \eta^2 = .28, p < .001$ ) and also within both the children ( $F(1, 52) = .1, \beta = -0.056, \eta^2 = .20, p < .001$ ) and adults ( $F(1, 49) = 9.1, \beta = -0.096, \eta^2 = .16, p < .001$ ) taken individually (see Figure 6a). Within the children, age was inversely associated with scene complexity, with an average of 0.35 fewer objects per scene for each additional year  $F(1, 52) = 12.6, \eta^2 = .19, p < .001$ . Aside from this difference, we also assess whether children's or adults' scenes bear the hallmarks of being driven by confirming or distinguishing between a small set of possible rules.

If participants do follow a control of variables, confirmatory, or iterative counterfactual approach, we would expect the scenes generated by participants to be more similar to the initial example or one of their own preceding scenes, than to a random scene or a scene drawn from a different learning problem. If they are rather maximising information with respect to a larger set of hypotheses, or exploring the data space efficiently, we would expect the opposite pattern of independence or anticorrelation. To explore this, we constructed a distance metric that we used to measure the feature-dissimilarity between any pair of scenes. The metric is based on 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. This involved  $z$ -scoring and combining a "minimal-edit set" of feature differences and incorporating a proportional cost for additional or omitted objects and scaling by the number of objects in the scenes. We provide a detailed procedure and example of how we computed these edit

**Figure 6**

(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 C for detailed procedure and similarity matrices separated by component.

distances and break them down into their separate components in the Appendix C. 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 primarily *between* rather than *within* subject for children's choices. Within subject but across trials, the average inter-scene dissimilarity for children was  $3.60 \pm .33$  similar to that for adults'  $3.65 \pm .22$ ,  $t(102) = .83, p = .4$ . Focusing more narrowly, within the scenes produced by an individual subject while learning about a single rule, we see a reversal of the aggregate pattern. That is, within a learning task, children's scenes are marginally *less* diverse on average than adults' (children:  $3.30 \pm 0.459$ ,

613 adults:  $3.44 \pm 0.33$ ,  $t(102) = 1.77$ ,  $p = 0.08$ , Figure 6b&c).

614 Figure 6c breaks down the within-trial scene dissimilarity by test position for the  
 615 two agegroups. Adults' scenes are clearly anchored to the initial example (right hand  
 616 facet)—shown by the dark shading in the top row indicating high similarity decreasing from  
 617 left to right for later tests—Adults' scenes also look sequentially self-similar—shown by the  
 618 relatively darker shading along the diagonal compared to the off-diagonal. In contrast,  
 619 children's similarity patterns look more uniform. However, for both adults and children,  
 620 the first self-generated scene is more similar to the initial example than any other scene.

## 621 Interim Discussion

622 In sum, in our experiment we found children were only moderately less able to guess  
 623 rules that fit the evidence than adults and there were only moderate differences in the  
 624 compatibility between children's and adults' rules and their subsequent generalizations.  
 625 Most striking was the fact that children appeared to overfit the evidence more, producing  
 626 more complex, perhaps more naïve, characterizations of the rule-following scenes than did  
 627 adults. This can be seen in the larger number of quantifiers and relations mentioned in  
 628 children's rules than in adults', essentially referring to more different objects and more  
 629 complex properties of the learning scenes that were actually irrelevant to their label. As  
 630 well as generating more complex concepts, children created more complex test scenes that  
 631 appeared to be more repetitive overall, yet also appeared to be varied less systematically  
 632 than adults'.

## 633 Model comparison

634 To explore the basis for the diversity of guesses and generalizations, and of the  
 635 differences between children and adults' learning, we now turn to model-based  
 636 characterization of the behavioral data. We focus first on the guesses, then the  
 637 generalizations, and finally the scene creation. We will assess whether participants guess  
 638 and generalization patterns are better captured by Bayesian inference over samples from an  
 639 expressive latent prior—Probabilistic Context Free Generation (PCFG)—or rather by the  
 640 partially bottom-up generation—Instance Driven Generation (IDG) limited to hypotheses  
 641 inspired by patterns in scenes (Bramley et al., 2018). We then assess whether new scenes  
 642 are better captured as independently generated—consistent with uncertainty-driven or  
 643 exploration-driven testing—or as adaptations of earlier scenes—consistent with  
 644 confirmatory or iterative contrastive testing.

645 To foreshadow, we find convergent evidence that both children's and adults' guesses  
 646 are better accounted for by Instance Driven Generation (IDG) of hypotheses than by an

647 approximately normative Probabilistic Context Free Grammar (PCFG) norm. We then  
 648 demonstrate that neither children’s nor adults’ generalizations can be explained by surface  
 649 similarity between rule-following and generalization probe scenes, but that they are well  
 650 predicted by the learners’ own symbolic guess. Finally, we show that almost all children’s  
 651 and adults’ scenes are more likely to have been created by making simplifications and edits  
 652 to either the previous or the initial scene—in line with hypothesis-driven confirmatory or  
 653 contrastive testing—rather than being generated independently from scratch—consistent  
 654 with uncertainty-driven or direct exploration of the data space.

## 655 **Guesses**

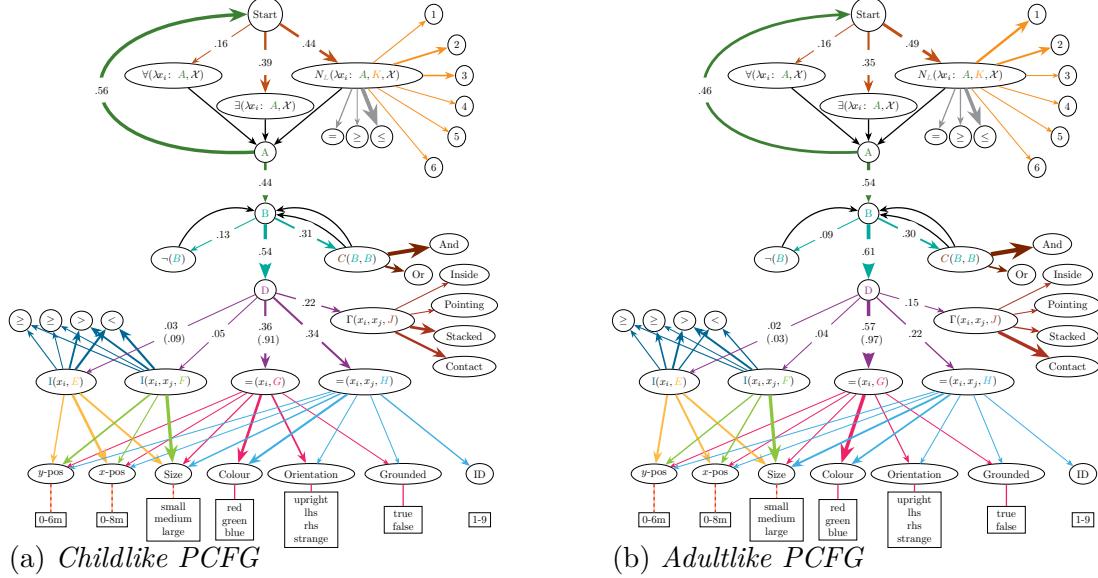
656 Participants produced a huge variety of guesses but despite this, their guesses were  
 657 consistent with the majority of their evidence. Children’s guesses were more complex and a  
 658 little less data-consistent on average than adults’. We now explore using PCFG and IDG  
 659 sampling to produce similar guesses.

660 We first assume a PCFG as a computational level framework and reverse engineer  
 661 what production weights it requires to generate the kinds of guesses we see adults and  
 662 children make. Next, we contrast the top-down first PCFG approach to rule generation  
 663 with our proposed data-inspired IDG, showing that the IDG does a better job of capturing  
 664 participants’ accuracy by problem type and agegroup and is also better able to produce the  
 665 specific guesses made by the participants.

### 666 ***Reverse engineering Childlike and Adultlike production weights***

667 Having encoded all the rule guesses from adults and children (in the section on *Rule*  
*668 complexity and constituents*), we created PCFG production weights that produce similar  
 669 guesses as adults and children. To do this, we worked back from the observed counts for  
 670 each rule element doing this separately for children’s and for adults’ guesses (see Appendix  
 671 A). Of course, the guesses are samples from a range of different participants’ posteriors,  
 672 since guesses were always based on some evidence. However, since this evidence differs  
 673 dramatically between trials and across the rules we considered and scenes participants  
 674 created, and since the structural elements of the grammar (booleans, quantifiers etc) are  
 675 not tightly tied to scene-specifics, this still provides a helpful elucidation of generation  
 676 differences behind child-like and adult-like guesses. A full set of fitted prior weights for  
 677 both adults and children are visualized in Figure 7. This analysis simply demonstrates that  
 678 a natural way to understand children’s guesses are as emanating from a less fine-tuned  
 679 generation mechanism adults’, with flatter, more entropic branching at 12 of the 14 forking  
 680 production steps we assumed in our PCFG model. Indeed probability distribution over

- 681 productions at each stage averaged  $1.28 \pm 0.50$  bits for children compared to  $1.03 \pm 0.59$   
 682 bits for adults,  $t(13) = 3.2, p = 0.007$ .



**Figure 7**

Visualization of (a) child-like and (b) adult-like PCFGs, reverse engineered to produce rules with empirical frequencies matched to children’s and adults’ guesses. 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.

### 683 Modeling accuracy by participant and rule

We now ask whether children’s or adults guesses depart from normative framework in being inspired by patterns present in particular learning scenes. To do this we compare participants patterns of accuracy to simulated approximately normative inference over a PCFG-generated sample and IDG hypothesis generation algorithms provided with the active learning data generated by the human participants. We generated a sample of 10,000 hypotheses based on uniform production weights  $\hat{H}_{\text{PCFG}_U}$ , and similarly for the IDG generated a sample based on uniform productions for each task  $\hat{H}_{\text{IDG}_U}^{p,t}$ . Additionally, for each participant  $p$ —and separately for each learning task  $t$  in the case of the IDG—we generated another 10,000 possible rules using age-consistent prior production weights derived above  $\hat{H}_{\text{PCFG}_H}^p$  and  $\hat{H}_{\text{IDG}_H}^{p,t}$  that have statistics matched to those in Figure 4a–f.<sup>7</sup> The PCFG samples act as an approximation to an infinite latent prior over rules  $P(h)$

<sup>7</sup> For these, we held out the subjects own guesses when setting the weights to avoid double dipping the data.

before seeing any data. The uniform-weight PCFG samples capture a generic inductive bias for simpler hypotheses while fitted held-out child- and adult-like weights additionally attempt to capture “learned” inductive biases common to the requisite age-group (but not specific to the participant). The IDG samples are additionally idiosyncratically constrained in the sense of only reflecting rules referring to features or relations actually present in at least one of the learning scenes. We split the IDG sample evenly across tests such that 1250 were “inspired” by each learning scene, necessarily repeating this procedure for each trial for each participant since each generates different evidence. In order to approximate a posterior over rules given self-generated learning scenes  $\mathbf{d}$ , we then weighted these samples by their likelihood of producing all eight scene labels  $l$  observed during the learning phase

$$P(h|l; \mathbf{d}) \propto P(l|h; \mathbf{d})P(h) \quad (1)$$

$$\approx P(l|h; \mathbf{d}) \sum_{\hat{h} \in \hat{H}} \mathbb{I}(h = \hat{h}) \quad (2)$$

and combined this with their prior weight—given by counting how often they appear in the prior sample, with indicator function  $\mathbb{I}(\cdot)$  denoting exact or semantic equivalence. To test for semantic equivalence, we computed predictions for the first 1000 participant-generated scenes for each rule and clustered together those that made identical predictions. We rounded positional features to one decimal place in evaluating rules to accommodate perceptual uncertainty. Concretely, we assumed the following likelihood function

$$P(l = 1|h; \mathbf{d}) \propto \exp(-b \times N_{\text{mispredictions}}) \quad (3)$$

embodying the idea that: the more learning scene labels a rule cannot explain, the less likely it is to have produced them. For a large  $b$ , the likelihood function approaches the true deterministic behavior of the rules. However, in our analyses we simply assume a  $b = 2$  to allow for some noise while maintaining computational tractability. This corresponds to a likelihood function that decays rapidly from  $\propto 1$  for rules that predict all 8 scenes’ labels, to  $\propto .13$  for a single misprediction, and  $\propto .02$  for 2 mispredictions, and so on.

To generate IDG predictions, we merged the production probabilities from the PCFG into the Instance Driven Generation procedure detailed in the Appendix A. For scenes that did not follow the rule we followed the same procedure as for scenes that did, but wrapped the rule in a negation. For example, observing a non-rule-following scene in which there are objects in contact might inspire the rule that “no cones are touching”.

The resulting model guess accuracy is shown in Table 2 and visualized in Figure 5a. We distinguish between two possible decision mechanisms: (1) Taking the *maximum*  $a$

703 *posteriori* (MAP) estimate from a large posterior sample (guessing in the event of ties),  
 704 which we take as closer to a normative ideal and (2) taking the accuracy of a single  
 705 posterior sample, which we take to be more consistent with the best-case-scenario output  
 706 of a process in which a given learner searches over hypotheses driven by a combination of  
 707 prior complexity and fit. Under all models, the MAP lines up with the correct hypothesis  
 708 more often than participants do (15–37% based on children’s active learning and 20–51%  
 709 based on adults’). For instance, under a uniform-weighted prior sample, the PCFG MAP is  
 710 correct on  $15\% \pm 35\%$  of children’s trials and  $20\% \pm 40\%$  of adults’ trials. Note that since  
 711 these simulations use the same prior sample, the small differences we see are due to the  
 712 different learning data generated by children and adults. However, accuracy improves  
 713 substantially and better reproduces the empirical child–adult accuracy difference when we  
 714 use samples based on reverse engineered weights that reproduce the qualitative properties  
 715 of other participants in the same agegroup (see Appendix A and Figure 7). For  
 716 age-appropriate prior samples, the PCFG guesses correctly on  $18\% \pm 38\%$  of children’s trials  
 717 and  $32\% \pm 46\%$  of adults’ trials. Using an age-inappropriate “flipped” prior sample (i.e.  
 718 child-like weights for adults and adult-like weights for children) obliterates this difference,  
 719 resulting in  $23\% \pm 42\%$  for children and  $22\% \pm 41\%$  for adults. We see a similar pattern for  
 720 the IDG algorithm, but higher accuracy across the board. The IDG achieves the best  
 721 accuracy on both children’s and adults’ trials, guessing over half of the hidden rules  
 722 correctly ( $51\% \pm 50\%$ ) in the case of adults’ trials. However, achieving this level requires  
 723 maximizing over the full sample, while we have argued that process level accounts are more  
 724 likely to yield behavior closer to posterior sampling (Table 2, right hand columns). Indeed  
 725 posterior samples provide a visually closer fit to the by-rule guess rates (Figure 5a).

726 To check what provides the better account of participants trial-by-trial accuracy  
 727 patterns we fit logistic mixed-effect regression models using each algorithm and prior  
 728 combination to predict each participant’s by-task probability of guessing correctly,  
 729 including random effects for both rule type and participant. The “Fit” columns of Table 2  
 730 shows the log likelihood for each of these models, revealing that participants’ correct  
 731 judgments were best predicted by posterior sampling under an IDG prior, with an  
 732 age-appropriate production weights (log likelihood = 211.5,  
 733  $\beta = 5.44 \pm 1.74$ ,  $Z = 5.99$ ,  $p < .001$ ) improving over a baseline fit of -234.3 for a model with  
 734 only intercept and random effects.

### 735 ***Modeling rule guess***

736 As a more direct test of the constructivist PCFG and IDG models’ ability to explain  
 737 participants’ free response guesses, we also attempted to estimate the probability of each

**Table 2**  
*Accuracy of Rule Guesses by Simulation Models*

| Algorithm  | Prior           | Accuracy MAP (%) |              |             | Accuracy Posterior Sample (%) |                |             |
|------------|-----------------|------------------|--------------|-------------|-------------------------------|----------------|-------------|
|            |                 | Children's data  | Adults' data | Fit         | Children's data               | Adults' data   | Fit         |
| PCFG       | Uniform         | 15±35            | 20±40        | -231        | 9±12                          | 12±14          | -227        |
| PCFG       | Agegroup        | 18±38            | 32±46        | -233        | 12±19                         | 20±25          | -228        |
| PCFG       | Flipped         | 23±42            | 22±41        | -235        | 16±21                         | 15±22          | -231        |
| IDG        | Uniform         | 27±44            | 39±48        | -228        | 9±12                          | 14±5           | -218        |
| <b>IDG</b> | <b>Agegroup</b> | <b>37±48</b>     | <b>51±50</b> | <b>-229</b> | <b>14 ± 16</b>                | <b>24 ± 22</b> | <b>-216</b> |
| IDG        | Flipped         | 26±44            | 52±50        | -234        | 14±20                         | 23±22          | -227        |

“Children” and “Adults” columns show the  $M \pm SD\%$  correct accuracy of the requisite algorithm based on the learning data from that agegroup. “Fit” shows the log likelihood for a logistic mixed-effects regression using this model to predict the probability of each participant guessing correctly on each trial.

738 approach generating exactly the participant’s encoded guess based on their active learning  
 739 data.

740 By definition, all 87% of trials in which participant gave an unambiguous rule, we  
 741 were able to encode in our concept grammar, so all have nonzero support under a PCFG  
 742 prior. Due to the stochasticity we assumed in our likelihood function, all possibilities also  
 743 nonzero have posterior probability, meaning they are guaranteed to appear in a sufficiently  
 744 large PCFG sample.<sup>8</sup> However, in practice it is impossible to cover an infinite space of  
 745 discrete possibilities with a finite set of samples, meaning there are a substantial number of  
 746 cases in which we did not generate the participants’ guess. The proportion of rules that  
 747 were generated at least once in 10,000 samples with agegroup fitted weights was highest for  
 748 the IDG with fitted weights (69% for children 76% for adults), decreasing to 49% and 62%  
 749 using uniform weights. This was still higher than for the PCFG which generated 42% for  
 750 children’s and 53% for adults’ guesses with the fitted prior weights and 45% for children’s  
 751 and 50% for adults’ rules from a uniform prior.

752 Table 3 details model fits to participants’ guesses. The IDG is again the stronger  
 753 hypothesis generation candidate, assigning higher probabilities on average to the rules that  
 754 participants provided. As expected, the variants of the PCFG and IDG with  
 755 agegroup-consistent production weights were better aligned with participants’ guesses than  
 756 variants with uniform (or mismatched) weights. However, all models produced adults’

---

<sup>8</sup> They would not necessarily appear in an infinitely large IDG sample because many of the more complex concepts are merely possible without being positively present. For example “there is a red and fewer than five small blues” is consistent with the Figure 1b but would never be generated by the IDG procedure inspired by these scenes.

**Table 3***Model Probability of Producing Participants' Exact Rule Guesses*

| Algorithm  | Prior           | Children         |           | Adults             |           |
|------------|-----------------|------------------|-----------|--------------------|-----------|
|            |                 | Mean (%)         | N best    | Mean (%)           | N best    |
| PCFG       | Uniform         | 3.3 ± 5.0        | 13        | 7.2 ± 7.2          | 10        |
| PCFG       | Agegroup        | 4.3 ± 7.4        | 13        | 12.5 ± 12.0        | 15        |
| IDG        | Uniform         | 3.4 ± 5.1        | 10        | 8.7 ± 8.6          | 2         |
| <b>IDG</b> | <b>Agegroup</b> | <b>4.5 ± 7.1</b> | <b>15</b> | <b>14.1 ± 13.6</b> | <b>22</b> |

Note: N best columns show the number of participants in each agegroup best fit by each model.

757 guesses with a much higher probability than children's guesses.

758 Figure 8a additionally visualizes participants' guesses in terms of their posterior  
 759 probability under PCFG and IDG sampling and compares this to what we would expect if  
 760 guesses are samples from the posterior (black line), the result of finding the maximum a  
 761 posteriori guess of the 10,000 considered hypotheses (dashed line) or else are simply  
 762 samples from the prior (dotted line). This visualization shows that, under all the models  
 763 we consider, adults' guesses are distributionally more consistent with posterior sampling  
 764 than posterior maximization, while children's appear somewhere between prior and  
 765 posterior sampling.

766 To better understand why we were not able to generate all of participants' guesses,  
 767 we also examined those frequently generated by the models and contrasted these with those  
 768 never generated under any of our model variants. Table 4 shows two examples of each for  
 769 children and adults and the full set is available in the [Online Repository](#). Unsurprisingly,  
 770 the participant guesses our models failed to generate tended to have more complex forms  
 771 and a concomitantly low generation probability. Assuming uniform weights, the syntax of  
 772 the children's guesses that we did generate had marginally higher log prior generation  
 773 probabilities Median (Inter-Quartile Range) -10.2 (5.0) than those we didn't were unable to  
 774 generate -13.9 (16.31) (Mood's median test,  $Z = 1.9, p = 0.053$ ). For adults this difference  
 775 was more pronounced -9.9 (5.0) compared to -14.9 (14.0) (Mood's median test,  
 776  $Z = 4.5, p = < .001$ ).<sup>9</sup> This examination revealed that one class of rules our participants  
 777 guessed but our models did not generate were those that could be expressed much concisely  
 778 with more powerful logical grammar. For example, we saw a number of cases of universal  
 779 quantification over feature values, such as "one of each color", mentioned in both a child

<sup>9</sup> Note that these prior generation probabilities are a lower bound on the chance of generating a particular semantic rule since many syntactic forms can express the same semantic content (Fränken et al., 2022). This captures why some relatively frequently generated semantic classes of guess nevertheless had a low probability for each specific syntactic expression .

**Table 4**  
*Example Guesses*

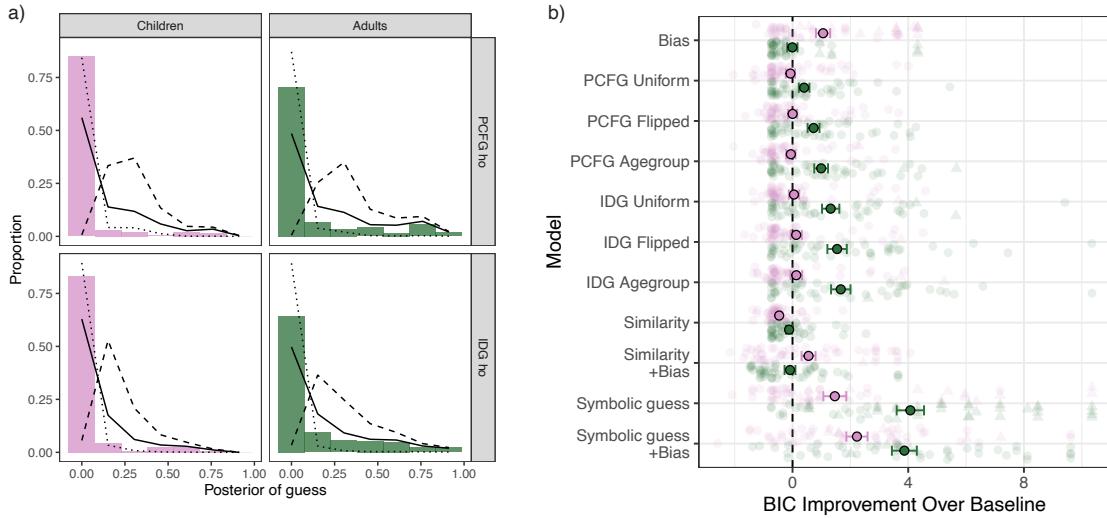
| Agegroup | Rule  | Example syntax  | log Prior<br>Uniform | log Prior<br>Age-<br>group | log(Likelihood)<br>N/10k |     |
|----------|---|---|----------------------|----------------------------|--------------------------|-----|
| Children | <i>“One is on top of the other”</i>   | $\exists(\lambda x_1 : \exists(\lambda x_2 : \Gamma(x_1, x_2, \text{stacked}), \mathcal{X}), \mathcal{X})$  | -9.5                 | -8.4                       | 0                        | 117 |
| Children | <i>“Only different colors”</i>  | $\forall(\lambda x_1 : \forall(\lambda x_2 : \vee(= (x_1, x_2, \text{ID}), \neg(= (x_1, x_2, \text{color}))), \mathcal{X}), \mathcal{X})$   | -9.8                 | -8.0                       | 0                        | 260 |
| Adults   | <i>“If there are multiple small blocks.”</i>  | $N_{\geq}(\lambda x_1 := (x_1, 1, \text{size}), 2, \mathcal{X})$  | -9.9                 | -19.6                      | 0                        | 609 |
| Adults   | <i>“There is at least one small green triangle.”</i>                                  | $\exists(\lambda x_1 : \wedge(= (x_1, \text{green}, \text{color}), = (x_1, 1, \text{size})), \mathcal{X})$  | -13.8                | -21.3                      | 0                        | 532 |
| Children | <i>“They have to be with all three different colors”</i>                              | $\exists(\lambda x_1 : \exists(\lambda x_2 : \exists(\lambda x_3 : \wedge(\wedge(= (x_1, \text{red}, \text{color}), = (x_2, \text{green}, \text{color})), = (x_3, \text{blue}, \text{color})), \mathcal{X}), \mathcal{X}), \mathcal{X})$                                    | -22.3                | -16.6                      | -2.0                     | 0   |
| Children | <i>“There has to be one small blue piece and there has to be more than one piece”</i> | $\exists(\lambda x_1 : N_{\geq}(\lambda x_2 : \wedge(= (x_1, 1, \text{size}), = (x_1, \text{blue}, \text{color})), 2, \mathcal{X}), \mathcal{X})$   | -12.5                | -11.3                      | 0                        | 0   |
| Adults   | <i>“When there is a cone from each color of the same size”</i>                        | $\exists(\lambda x_1 : \exists(\lambda x_2 : \exists(\lambda x_3 : \wedge(\wedge(\wedge(= (x_1, \text{red}, \text{color}), = (x_2, \text{green}, \text{color})), = (x_3, \text{blue}, \text{color})), = (x_1, x_2, \text{size})), \mathcal{X}), \mathcal{X}), \mathcal{X})$ | -20.5                | -11.11                     | -2.0                     | 0   |
| Adults   | <i>“one piece has to be leaning on another”</i>                                       | $\exists(\lambda x_1 : \exists(\lambda x_2 : \wedge(\Gamma(x_1, x_2, \text{contact}), \neg(= (x_2, \text{upright}, \text{orientation}))), \mathcal{X}), \mathcal{X})$   | -18.5                | -21.3                      | -3.9                     | 0   |

Note N/10k shows how many times we generated this rule in 10,000 samples assuming agegroup-specific weights and counting any semantically equivalent expressions.

and an adult guess in Table 4. This kind of rule can be expressed parsimoniously in second order logic with a single universal quantifier over color properties while in our grammar it required a separate quantification for each color. The fact that children produced about as many apparently higher-order-logic rules as adults seems to suggest that the PCFG we assumed, despite its ostensibly complex structure, is still a simplification of the basis from which children constructed their ideas (cf. Piantadosi et al., 2016).

## 786 Generalizations

We next examine our models’ ability to account for participant’s generalization performance. As with the guesses, we first examine patterns of accuracy by comparing participants to simulated constructivist PCFG and IDG learner benchmarks before fitting a range of models to the specific generalizations participants made.

**Figure 8**

a) Posterior probability of participants' guesses under PCFG and IDG samples. Full black line compares with posterior samples, dashed line with selection of the posterior maximum a posteriori hypothesis (or sampling from them if there are more than one), dotted line compares with samples from the prior. b) Individual generalization model fits showing BIC improvement over baseline per trial (higher is better). Opaque points show mean $\pm$ SE, faint points show individual fits, with triangles used to mark where the model is the best fit (of all 18 tested) for that participant.

### 791 ***Modeling generalization accuracy***

To do this, we use their requisite predictive distributions to model labelling generalizations  $\mathbf{l}^*$  to the set of test scenes  $\mathbf{d}^*$

$$P(\mathbf{l}^*|\mathbf{l}; \mathbf{d}, \mathbf{d}^*) = \int_H P(\mathbf{l}^*|H; \mathbf{d}^*)P(H|\mathbf{l}; \mathbf{d}) dH \quad (4)$$

$$\approx \sum_{h \in \hat{H}} P(\mathbf{l}^*|h; \mathbf{d}^*)P(h|\mathbf{l}; \mathbf{d}) \quad (5)$$

792 Provided with the active learning data generated by the human participants, both  
 793 performed in the human range at generalization. As with predicting the guesses, taking the  
 794 marginally most likely generalization labels over a posterior weighted sample of  
 795 agegroup-appropriate IDG prior productions performed best overall and reproduced the  
 796 difference between children's and adults' generalization accuracies ( $68.8 \pm 20.1\%$  and  
 797  $74.2\% \pm 21.7\%$ ). The uniform-production IDG still performed slightly better than the  
 798 PCFG, generalizing at  $65.2\% \pm 19.3\%$  from children's active learning data and  
 799  $69.0\% \pm 21.0\%$  from adults'. Using agegroup-appropriate priors, the PCFG also reproduces  
 800 the empirical difference between children's and adults' accuracy:  $62.8 \pm 19.8\%$  for children's

801 trials and  $68.8 \pm 20.9\%$  for adults' trials. Using the PCFG with uniform production weights  
802 yielded accuracies of  $61.4\% \pm 19.6\%$  for children's and  $63.5\% \pm 20\%$  for adults' data.

803 The stronger generalizations of the IDG compared to the PCFG replicates the  
804 findings of Bramley et al. (2018) and extends this to children as well as adults. Intuitively,  
805 this is because the bottom-up inspiration mechanism ties the hypotheses generated to  
806 features of the learning cases, effectively narrowing in on plausible hypotheses more  
807 efficiently. More broadly, these simulation results underscore the inherent difficulty of this  
808 task in particular and open-ended inductive inference in general. The PCFG and IDG were  
809 not statistically better or worse than participants at any rule inference after Bonferroni  
810 correction with the exception that the IDG outperformed children on rule 4  
811  $t(96) = 4.7, p < .0001$ . Thus strikingly, even in this "small world" with known and fully  
812 observed features, and even allowing simulations to sample and maximize over implausibly  
813 large numbers of hypotheses, we could not robustly outperform human adults in this task.  
814 This also reveals that building in human inductive biases boosts generalization  
815 performance (cf Lake et al., 2017) and the idea that adults' have formed stronger inductive  
816 biases than children goes some way to explain differences in how they generalize.

817 A complicating factor is that children generated different learning data to adults.  
818 However, our PCFG and IDG simulations suggest exposure to different data cannot explain  
819 most of the accuracy differences between children and adults. Using identical production  
820 weights and the scenes generated by adults and children led to only small differences in  
821 accuracy for the PCFG and moderate for the IDG, while using a "flatter" set of productions  
822 fit to match childlike rules, and a more "peaked" set fit to adults' rules, better reproduces  
823 the accuracy differences. We take this to suggest hypothesis construction differences drive  
824 a large portion of the differences in children's and adult's inductive inferences.

### 825 *Modeling specific generalizations*

826 A standard benchmark for models of concept learning is a fit with participants'  
827 generalizations to new exemplars. Thus, we compared a range of models' ability to account  
828 for participant's specific generalizations. The set of models we consider allows us to test  
829 our core claims that children's and adults' induced representations are symbolic and  
830 compositional, as opposed to statistical and similarity-based.

831 We fit a total of 18 models to the generalization data. All models had between 0  
832 and 2 parameters. For each model, we fit the parameter(s) by maximizing the model's  
833 likelihood of producing the participant data, using R's `optim` function. We compared  
834 models using the Bayesian Information Criterion (Schwarz, 1978) to accommodate their  
835 different numbers of fitted parameters.

836 The models we fit were:

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

840 **2. Bias.** Acts a stronger baseline by allowing participants to have an overall bias  
 841 toward or against selecting generalization scenes as rule following. For this model,  $b$   
 842 = 1 if >50% of generalizations predict the scene is rule following and 0 otherwise.  
 843 The model is fit using a mixture parameter  $\lambda$  to mix this modal prediction with the  
 844 baseline prediction of .5  $P(\text{choice}) = \lambda b + (1 - \lambda).5$ .

845 **3-8. PCFG {Uniform, Flipped, Agegroup} {No Bias, Bias}.** These models  
 846 base their generalizations on the marginal likelihood that each generalization scene is  
 847 rule following under the Probabilistic Context Free Generation (PCFG) posterior  
 848  $r = P_{\text{PCFG}}(\mathbf{l}^* | \mathbf{l}; \mathbf{d}, \mathbf{d}^*)$ . “Uniform” uses a prior with uniform production weights.  
 849 “Flipped” uses a prior generated with mismatched weights — that is, adultlike  
 850 weights for children’s generalizations and childlike weights for adults’ generalizations.  
 851 “Agegroup” uses a sample based on weights derived from other participants in the  
 852 same agegroup holding out the participants’ own guesses. In each case, these  
 853 predictions are then softmaxed using  $P(\text{choice}) = \frac{e^{r/\tau}}{\sum_{r \in R} e^{r/\tau}}$ , with temperature  
 854 parameter  $\tau \in (0, \infty)$  (Luce, 1959) optimized to maximize model likelihood. Large  
 855 positive  $\tau$  indicates random selection.  $\tau \rightarrow 0$  indicates hard maximization. Variants  
 856 with a bias term also mix this prediction with the subject’s modal response  $b$  as in

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

857 **9-14. IDG {Uniform, Flipped, Agegroup} {No Bias, Bias}.** These models use  
 858 the marginal likelihood of each generalization scene as rule following under the  
 859 Instance Driven Generation based posteriors with variants as with the PCFG variants  
 860 and again fit with softmax parameter  $\tau \in (0, \infty)$ .

861 **15-16. Similarity {No Bias, Bias}.** Inspired by Tversky’s statistical and  
 862 similarity based *contrast model of categorization* (cf., Tversky, 1977), we used the  
 863 inter-scene similarity between each generalization scene and each training scene to  
 864 compute the relative average similarity of each generalization case to the  
 865 rule-following vs. the not rule-following training scenes. Similarities were computed  
 866 using the same procedure used in the Active Learning section of the Results and

detailed in Appendix C. We computed the mean difference between rule-following and not-rule following similarities as a  $\Delta\text{Similarity}$  score for each participant  $\times$  trial  $\times$  item combination. Positive scores mean generalization item has a greater feature similarity to the rule following learning scenes than the not rule-following learning scenes. Negative scores mean the reverse. To convert these into choice probabilities, we take a logistic function of these scores  $r = \frac{e^{\Delta\text{Similarity}}}{e^{\Delta\text{Similarity}} + 1}$  and again fit these  $r$  values to maximize the likelihood of participants' choices using a softmax function with inverse temperature parameter  $\tau \in (0, \infty)$ . Intuitively, this model provides a non-symbolic alternative account of generalization behavior.

**17-18. Symbolic Guess {No Bias, Bias}.** This model takes participants' free guess of the hidden rule, coded in lambda abstraction, and uses these directly to generate a prediction vector  $r \in R : \{\text{rule-following}=1, \text{not rule-following}=0\}$  for each scene. For trials in which the participant does not provide an unambiguous rule, the model assigns a .5 likelihood to each generalization choice. These were again fit with a softmax parameter  $\tau \in (0, \infty)$ .

A good fit for *Symbolic Guess* would support our core claim that participants inductive generalizations are directly driven by their constructed symbolic ideas. Meanwhile, a better fit for *Similarity* would suggest that generalizations are rather based on sub-symbolic feature similarity, with participants guesses relegated to a supporting role as rough symbolic re-descriptions of an ultimately sub-symbolic representation (e.g., Dennett, 1991; Johansson, Hall, & Sikström, 2008). To the extent that our constructivist simulations reflect participants' inductive inference mechanisms we expect the end-to-end PCFG and IDG models to also capture generalization patterns even though they are blind to the individual participants' explicit guesses. This also acts as a sanity check for our approach for any readers skeptical about the validity of self-report data.

We fit all models to the children's and adults' data, and then separately to each individual participant. The full table of model fits is presented in the Appendix (Table A-3). Individual level results are highlighted in Figure 8b. At the individual level, the PCFG+Bias and IDG+Bias models performed no better than the unbiased PCFG or IDG models, thus we omit these from Figure 8b for simplicity.

In line with our core hypothesis, *Symbolic guess + Bias* is the best fitting model of both children's and adults' generalizations outperforming all the models we considered based just on only the learning data. For children's generalizations taken together, *Symbolic guess + Bias* has BIC 2149, improving 490 over Baseline with bias term mixture weight of  $\lambda = .26$  and choice temperature parameter  $\tau = 0.80$ . For adults, this is BIC 1776

902 with a larger BIC improvement of 996 over Baseline, with a  $\lambda = 0.08$  indicating less bias  
 903 and temperature  $\tau = 0.50$  indicating tighter alignment with the guessed-rule's predictions.  
 904 Probing this bias, we see children undergeneralized substantially on average, selecting just  
 905  $2.75 \pm 1.42/8$  scenes compared to adults'  $3.42 \pm 1.03/8$  (unknown to the participants, there  
 906 were always 4 rule following generalization scenes). Focusing on individual fits, the picture  
 907 is mixed for children's generalizations, with 16/50 best fit by the *Bias* only model, followed  
 908 by 15 by the *Symbolic guess* model, 9 by the *Symbolic Guess + Bias* model and a further 7  
 909 by the fully random *Baseline*. No other model best fit more than 2 children. For adults,  
 910 32/52 were best fit by *Symbolic guess*, 6 by *Bias*, 4 by *Symbolic guess + Bias* and no other  
 911 model best fit more than 2 participants.

912 Overall, children's generalizations were much harder to predict than adults' with  
 913 end-to-end constructivist accounts of their generalizations performing close to *Baseline*.  
 914 This is partly to be expected since our child-like construction weights inherently produce a  
 915 very diverse set of guesses and correspondingly diffuse set of generalization predictions.  
 916 However, conditioning on Children's symbolic guesses we were able to predict their  
 917 generalizations far better than by *Similarity*, *Bias* or any other model we considered.  
 918 Adults' generalizations seem more straightforwardly driven by their symbolic guesses, with  
 919 better individual fits on average using their guess directly without adjusting by any bias  
 920 toward or against predicting scenes to be rule-following. This makes sense: with a clear  
 921 hypothesis in mind, there is little rationale to select more or fewer than the generalization  
 922 scenes consistent with that rule.

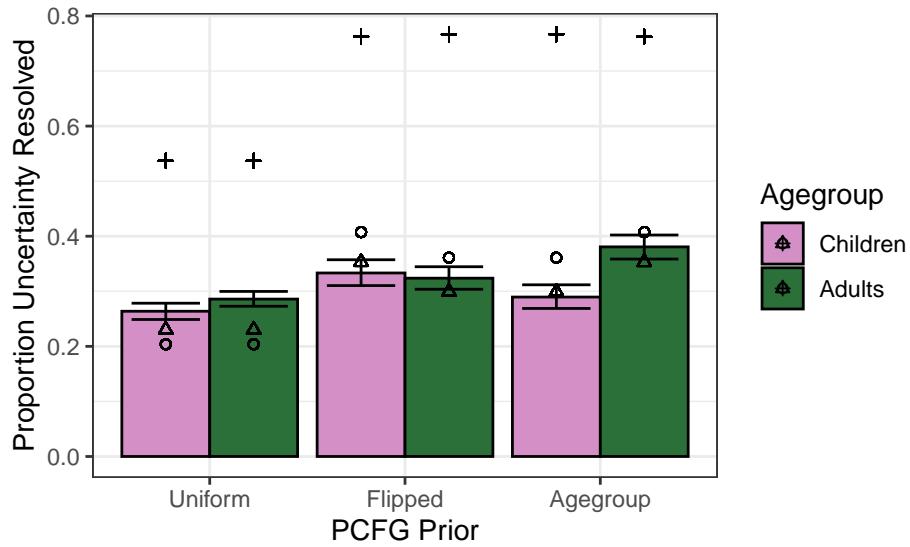
923 As with the free rule guesses, the IDG was robustly more aligned with participants'  
 924 generalizations than the PCFG, particularly for adults, and particularly when using  
 925 agegroup-appropriate weights rather than Uniform or age-inappropriate Flipped  
 926 production weights. Thus, this model comparison also supports the idea that participants  
 927 were inspired by patterns present in the learning data, such as the objects and relations in  
 928 the initial positive example. However, this does not appear to be a developmental  
 929 difference per se, with both children's and adults' judgments better accounted for by the  
 930 IDG than our PCFG algorithm across all analyses.

931 These results support a key aspect of the constructivist framework, participant's  
 932 idiosyncratic symbolic guesses seem to do the work in driving generalizations, rather than  
 933 these being driven by family resemblance in the features of the scenes. The constructivist  
 934 account anticipates that generalization patterns are dependent on what concept the learner  
 935 has arrived at by the end of learning, and our end-to-end models of this process  
 936 demonstrate the sheer breadth of concepts that learners can reasonably end up with in this  
 937 task.

938 **Scene generation**

939 We finally turn to participants' scene generation. We compare participants  
 940 generated scenes to several benchmarks before comparing a set of models of scene  
 941 generation to test the idea that participants adapted earlier scenes to isolate and test the  
 942 role of features mentioned in their hypotheses.

943 ***Comparison with information norms***



**Figure 9**

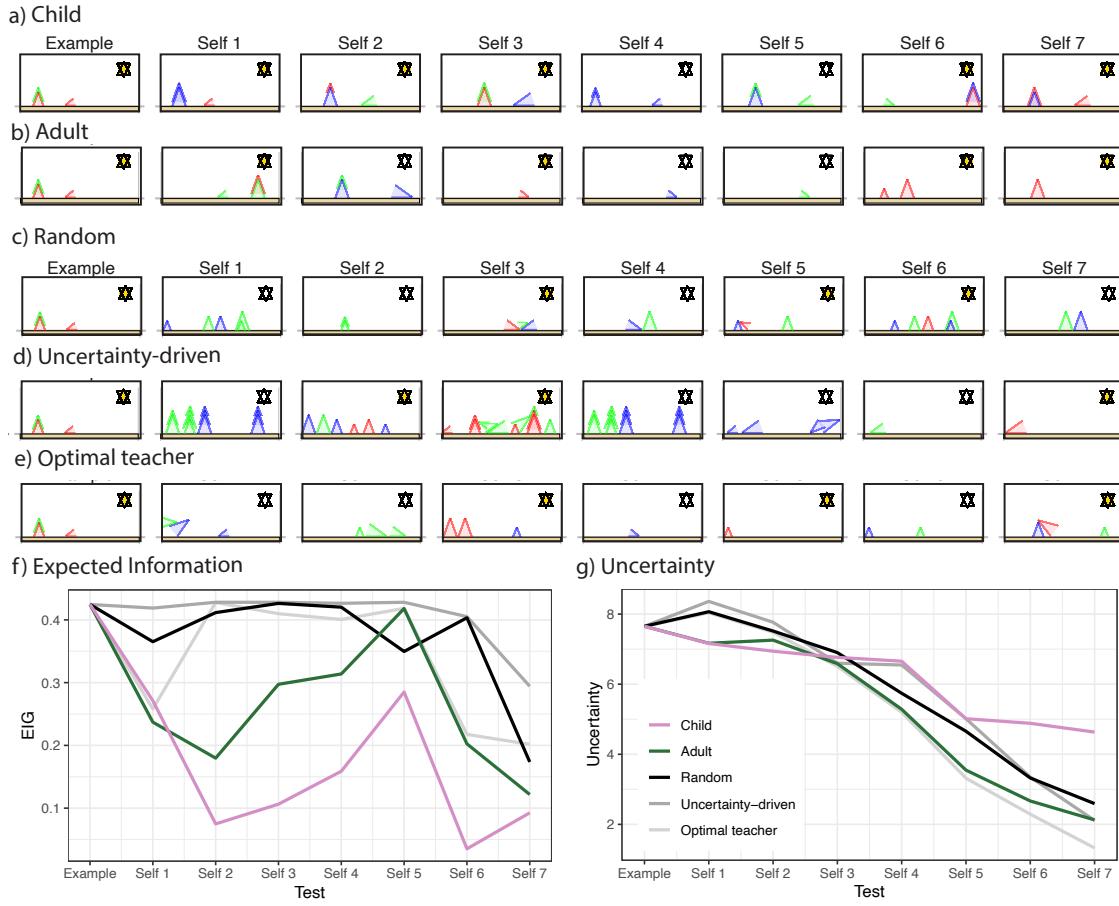
*Uncertainty reduction under different priors. Triangles = random scene selection. Circles = greedy expected information maximizing scene selection. "+" symbols = Ideal teaching scenes.*

944 According to an information gain analysis, children's and adults' scene generation  
 945 result in some differences in the quality of the total evidence generated. For example, using  
 946 the unweighted PCFG sample, prior entropy is 7.74 bits and children's evidence produces  
 947 an information gain (reduction in uncertainty) of  $1.93 \pm 0.45$  bits while adults' data average  
 948 an information gain of  $2.11 \pm 0.38$  bits  $t(102) = 2.12, p = 0.035$  (see Figure 9). Relative to  
 949 the agegroup-fitted PCFG priors, the difference in information gains is rather larger, with  
 950 children's scenes leading to information gain at  $2.28 \pm 0.66$  bits (prior entropy  $7.87 \pm 0.05$ ),  
 951 and adults' at  $2.96 \pm 0.64$  (prior entropy  $7.77 \pm 0.04$ )  $t(102) = 5.3, p < .0001$ . Under the  
 952 flipped priors—that is, taking the adultlike PCFG prior for children and childlike PCFG  
 953 prior for adults—children's tests look more informative than under their own prior,  
 954 generating  $2.58 \pm 0.68$  bits, and adults' tests slightly less informative than under their own  
 955 prior  $2.55 \pm 0.57$  bits, eliminating the statistical difference  $t(102) = 0.24, p = 0.81$ . On the

956 face of it, this is evidence against the idea that children’s more elaborate hypothesis  
 957 generation and concomitantly flatter construction weights are driving them rationally  
 958 toward more elaborate testing choices. However, as we noted information-theoretic  
 959 analyses as limited in what can reveal. It is predicated on an implausibly complete  
 960 representation of uncertainty that we approximated by using a large sample of prior  
 961 hypotheses, while we have characterized constructivist learning as driven by more focal  
 962 testing of a handful of similar possibilities.

963 We also compared participants against three scene selection benchmarks. In  
 964 Figure 9, black triangles show the reduction in uncertainty resulting from supplementing  
 965 the initial example with 7 scenes selected at random from among participant  
 966 generated scenes. Circles show the result of repeatedly selecting from a sample of 1000 of  
 967 the participant-generated scenes, greedily selecting whichever one maximizes the expected  
 968 information gain with respect to the prior at that test. Plus symbols show the reduction in  
 969 uncertainty resulting from observing scenes selected by an ideal teacher—i.e. the seven  
 970 scenes that, in combination with the initial example, best reveal the true concept. One  
 971 striking feature of these benchmarks is the low performance of the uncertainty-driven norm  
 972 under all PCFG priors. Expected information gain slightly outperforms participants and  
 973 random selection assuming the agegroup priors, but is actually worse than random scene  
 974 selection under a flat uniform prior sample. This poor performance stems from the fact  
 975 that the prior space of hypotheses is just so large and symmetric, making most scenes  
 976 similarly informative at first. Furthermore, a large class of PCFG hypotheses predict that  
 977 all possible scenes will be rule following, or that all possible scenes will be non-rule  
 978 following. These hypotheses are incorrect and rarely entertained by participants, yet have  
 979 an outsized effect on the greedy selection of scenes that maximize expected information  
 980 gain. Scenes selected to maximally convey each concept are far more informative,  
 981 highlighting gulf between self-teaching and optimal teaching in inductive settings.

982 Figure 10 compares an example scene sequence selected by a child and a child  
 983 against a random selection from all participant scenes, uncertainty-driven selection and  
 984 those selected to maximally convey the concept. This visual comparison highlights how  
 985 human scene selection involves recognizable repetition and patterning that look quite  
 986 unlike random and uncertainty-driven selection. In particular, several of the scenes selected  
 987 to minimize expected uncertainty are very complex compared to participants’ selections.  
 988 Theoretically uncertainty driven scenes do an excellent job of dividing the hypothesis  
 989 space, shown by their ceiling-level EIG (Figure 10f). However, since the target rule in this  
 990 case turns out to be a simple, this sophistication does not benefit the uncertainty-driven  
 991 learner overall (Figure 10g).

**Figure 10**

Example sequences for the “There is a red” problem. a) A child’s scenes b) An adult’s scenes c) Random selection from all participant generated scenes d) Uncertainty driven selection from all participant scenes e) Optimal scene selection for communicating the concept. f) Expected Information Gain and g) achieved uncertainty reduction for sequences in a–e.

## 992 *Models of scene selection*

993 We hypothesized participants might adopt incremental hypothesis-driven testing  
 994 strategies to deal with the challenges of the inductive setting. We suggested this might  
 995 involve testing nearby confirmatory generalizations of a focal hypothesis (Klayman & Ha,  
 996 1989), or contrasting nearby variants to this hypothesis (Oaksford & Chater, 1994). In  
 997 either case, we argued this would result in patterns of similarity (retention of rule-critical  
 998 elements and creation of minimal contrast pairs) and simplification (removal of non-rule  
 999 critical elements) quite distinct from the predictions of information-driven or  
 1000 uncertainty-driven testing. We indeed observed anchoring within learning problems. In  
 1001 particular, participants scenes appeared to be anchored both persistently to the initial

positive example and sequentially (Figure 6c). We here operationalize this by creating a family of scene adaptation models that assume learners create new scenes by mutating either the initial positive example, or their own previous scene. We compare these against baselines that rather assume learners generate each new scene from scratch. Concretely, the models we fit were:

1. **Generate {Uniform}**: Adds a random number of objects to each scene. Uniform assumes each object uniformly selected features (color, size, orientation and groundedness)<sup>10</sup>. This model has zero fitted parameters so acts as an overall baseline. Otherwise with this and all subsequent models we assumed each feature was sampled from its mean prevalence to act as a stronger baseline.
2. **Generate Simple**: Adds a number objects to each scene drawn from an exponential distribution (truncated to the maximum allowable number of objects) with fitted rate parameter  $\lambda$ , selecting the features of these objects at random. This models a tendency to create simple scenes containing fewer objects, with the mean number of objects per generated scene given by  $\frac{1}{\lambda}$ .
3. **Adapt Initial {Simple}**: Assumes the learner creates each new scene by adapting the initial scene. Concretely, we assume the learner samples either the same number of objects as in the initial scene with probability  $\eta$ , or a random number with probability  $1 - \eta$ . The objects in new scene are assumed to be a mixture of the features of the matching object in the initial scene (replicating the original feature with probability  $\eta$ ) or selected randomly from their support (with probability  $1 - \eta$ ). We marginalize over all possible object mappings between scene  $i$  and  $j$ .  $\eta = 1$  corresponds to perfectly reliable copying of the number and nature while  $\eta = 0$  denotes always resampling the feature. The simple variant assumes the number of objects in the scene, if not drawn from the inspiration scene, is drawn from an exponential distribution with parameter  $\lambda$  as above.
4. **Adapt Previous {Simple}**: This model works as above but uses the preceding scene rather than the initial scene as its starting point.
5. **Adapt Mixed {Simple}**: This model simply mixes the predictions of Adapt Initial and Adapt Previous to capture the behavior of a learner who sometimes adapts the initial scene (with probability  $\theta$ ) or by their own preceding scene with probability  $(1 - \theta)$ .

---

<sup>10</sup> We do not attempt to predict the relational features or absolute positions in this analysis.

**Table 5**  
*Models of Scene Generation*

| Model                 | Children  |    |                 | $\lambda$     | $\eta$        | $\theta$ |
|-----------------------|-----------|----|-----------------|---------------|---------------|----------|
|                       | BIC/scene | N  | Best            |               |               |          |
| Generate Uniform      | 40.2      | 0  |                 |               |               |          |
| Generate              | 34.9      | 0  |                 |               |               |          |
| Generate Simple       | 30.7      | 0  | $0.34 \pm 0.1$  |               |               |          |
| Adapt Initial         | 30.4      | 2  |                 | $.29 \pm .19$ |               |          |
| Adapt Previous        | 30.1      | 8  |                 | $.25 \pm .18$ |               |          |
| Adapt Mixed           | 30.0      | 1  |                 | $.27 \pm .19$ | $.40 \pm .29$ |          |
| Adapt Initial Simple  | 29.3      | 7  | $0.33 \pm 0.11$ | $.34 \pm .16$ |               |          |
| Adapt Previous Simple | 29.0      | 10 | $0.34 \pm 0.13$ | $.31 \pm .17$ |               |          |
| Adapt Mixed Simple    | 28.7      | 26 | $0.34 \pm 0.12$ | $.33 \pm .17$ | $.40 \pm .24$ |          |
| Adults                |           |    |                 |               |               |          |
| Model                 | BIC/scene | N  | Best            | $\lambda$     | $\eta$        | $\theta$ |
| Generate Uniform      | 32.8      | 0  |                 |               |               |          |
| Generate              | 27.8      | 0  |                 |               |               |          |
| Generate Simple       | 23.1      | 0  | $0.50 \pm 0.18$ |               |               |          |
| Adapt Initial         | 23.6      | 0  |                 | $.23 \pm .14$ |               |          |
| Adapt Previous        | 23.4      | 1  |                 | $.21 \pm .13$ |               |          |
| Adapt Mixed           | 23.3      | 1  |                 | $.21 \pm .13$ | $.35 \pm .26$ |          |
| Adapt Initial Simple  | 22.4      | 5  | $0.50 \pm 0.20$ | $.29 \pm .12$ |               |          |
| Adapt Previous Simple | 21.9      | 24 | $0.54 \pm 0.30$ | $.23 \pm .13$ |               |          |
| Adapt Mixed Simple    | 21.8      | 19 | $0.54 \pm 0.27$ | $.24 \pm .13$ | $.32 \pm .25$ |          |

Note: BIC/scene shows the fit of the model at the agegroup level divided by the number of scenes for easier comparison.  $\lambda$  (simplicity),  $\eta$  (fidelity) and  $\theta$  (mixture) show  $M \pm SD$  of best fitting model parameters variant across subjects.

We fit the models to each agegroup, and separately every individual participant (see Appendix B for details). Table 5 shows the resulting ageregroup-level BICs the number of individuals best fit by each model and the spread of parameter values for each. Adapt Mixed Simple was the best model for both agegroups overall and the best model for 48% of children and 38% of adults. No participant was better fit by Generate or Generate Simple, capturing that every single participant exhibited some degree of positive anchoring on the number or nature of the earlier scenes. 80% of children and 96% of adults additionally showed an additional preference for simple scenes. Almost half of adults (48%) were best characterized as adapting the previous scene than repeatedly adapting the initial scene or a mixture of both while this was only true for 19% of children. Fitted simplicity rate  $\lambda$  was larger for adults ( $\approx 0.5$ ) than children ( $\approx 0.3$ ) capturing their stronger tendency to create

1045 scenes with fewer objects. Fidelity of copying features of inspiration scenes  $\eta$  was similar  
1046 for children and adults ( $\approx .3$ ). Note that this is an underestimate due to the need to  
1047 marginalize over many possible object-object mappings and two potential inspiration  
1048 scenes. Mixture parameter  $\theta$  was below .5 on average for both children and adults  
1049 suggesting dominance of the initial scene over the previous scene.

1050 In sum, this model comparison supports the idea that learners adapted their earlier  
1051 tests often retaining the same number of objects and tending to keep many of the same  
1052 features. Adults were more likely than children to reduce the number of objects and had  
1053 more tendency to adapt sequentially, gradually traveling further away from the initial  
1054 example.

## 1055 General Discussion

1056 In this paper, we explored children and adults' active hypothesis generation and  
1057 inductive inference in an interactive task where the space of possibilities and actions is  
1058 compositional, open and practically unbounded. Our results are rich and nuanced but  
1059 broadly we found that:

- 1060 1. Children's guesses and tests were more complex than those of adults.
- 1061 2. We could synthesize the diversity and distribution of children and adults' guesses  
1062 with a constructivist—symbolic, generative—inference framework, reproducing both  
1063 their sporadic correct guesses but also capturing the spread of their incorrect ideas  
1064 and offering a framework for modeling differences between children's and adults'  
1065 inductive inference.
- 1066 3. Children's guesses reflected less fine-tuned construction mechanisms than adults',  
1067 producing more diversity but were consequently less predictable.
- 1068 4. Both children's and adults' hypothesis generation appeared data-inspired, shown by  
1069 better fit throughout our model-based analyses by our Instance Driven Generation  
1070 account—inspired by patterns in the learning scenes—over our approximately  
1071 normative (PCFG) account—that generated hypotheses a priori and weighted them  
1072 with the evidence.
- 1073 5. The logical form of both children and adults' symbolic guesses predicted their  
1074 generalizations to new scenes far better than feature similarity.
- 1075 6. Both children and adults scenes generation seemed to involve modifying previous  
1076 scenes, with adults doing so more systematically and with more tendency to simplify

1077 them.

1078 We now discuss these results more broadly, first highlighting some limitations, then  
1079 expanding on what we see as the implications of this work for theories of concepts and of  
1080 development and finally pointing to some future directions.

1081 **Limitations**

1082 ***Experimental Control***

1083 While this task and new dataset provide an exceptionally rich window on inductive  
1084 inference, some of what is gained in open-endedness is lost in experimental control. There  
1085 is considerable residual ambiguity about the extent that differences in active learning  
1086 shaped differences in hypothesis generation and visa versa. One way to try and partial this  
1087 out could be to run more experiments that fix the evidence and probe the hypotheses  
1088 generated, or that fix the hypotheses in play and probe what evidence is sought. However,  
1089 we have argued that such constrained tasks run the risk of short-circuiting natural  
1090 cognition: Learners may struggle to test hypotheses they did not conceive themselves, and  
1091 are known to struggle to use data they have not generated to evaluate their hypotheses  
1092 (Markant & Gureckis, 2014; Sobel & Kushnir, 2006). Sole focus scenarios fix one or other  
1093 aspect of the the inductive inference loop may provide a misleading perspective on  
1094 end-to-end active inference in the wild. We feel that our open ended task provides a  
1095 valuable complementary perspective. In future work hope, we plan to elicit more  
1096 fine-grained online measures of learners' thought process—e.g. asking them to list their  
1097 hypotheses after each guess or describe how they construct test scenes. This would support  
1098 comparison of process-level accounts of both hypothesis adaptation and active search and  
1099 allow identification of individual differences.

1100 ***Theoretical Expressivity***

1101 There are many ways we could have set up the primitives, parameters and  
1102 productions of our PCFG and IDG models. This makes for a dangerously expressive set of  
1103 theories of cognition. We do not claim to have explored this space exhaustively here but  
1104 rather that our modeling lends support to the idea that some symbolic and compositional  
1105 process drives children and adults' active inductive inferences about the world. That is, we  
1106 can explain the variability and productivity of human hypothesis belief formation in  
1107 symbolic terms. Identifying the computational primitives of thought may not be a realistic  
1108 empirical goal since a feature of constructivist accounts is their flexibility. Learners can  
1109 grow their concept grammar over time, caching new primitives that prove useful

1110 (Piantadosi, 2021). Moreover, it is well known many different symbol systems can mimic  
1111 one another (Turing, 1937), meaning that expressivity alone cannot distinguish between  
1112 them. Since, we expect different learners to take different paths in an inherently stochastic  
1113 learning trajectory, this limits universal claims about representational content.

1114 ***Feature selection***

1115 We assumed our scenes had directly observable features and cued these to  
1116 participants in our instructions. However, a number of recent models in machine learning  
1117 combine neural network methods for feature extraction with compositional engines for  
1118 symbolic inference, creating hybrid systems that can learn rules and solve problems from  
1119 raw inputs like natural images (cf. Nye, Solar-Lezama, Tenenbaum, & Lake, 2020; Valkov,  
1120 Chaudhari, Srivastava, Sutton, & Chaudhuri, 2018). We see these approaches as having  
1121 promise to bridge the gap between subsymbolic and symbolic cognitive processing.

1122 ***Elicitation differences between children and adults***

1123 One potential concern is that the complexity of children's guesses relative to adults  
1124 stems partly from their being collected verbally and in the presence of an experimenter  
1125 rather than typed during an online experiment. Speaking carries different cognitive  
1126 demands than typing and may lead to children simply responding in a more verbose way  
1127 than adults. While we cannot rule this out, we do not think this is a major concern.  
1128 Adults were well compensated for accuracy, meaning their motivation was primarily to be  
1129 correct rather than brief. The semantic content of both children's and adults' rules were  
1130 extracted through our coding of them into lambda calculus meaning that surface  
1131 differences in concise expression can be separated from logical complexity. Furthermore  
1132 children's guesses were not the only thing that was more elaborate about their behavior.  
1133 They were also more elaborate in their active testing choices, producing more complex  
1134 scenes despite having to create these in the same manner as adults. Since the testing  
1135 interface was reset on each trial, this complexity took more effort, with children's scenes  
1136 requiring substantially more clicks and more time to produce than adults'.

1137 ***Use of verbal protocols***

1138 Another worry about our use of free responses is that they rely on a capacity for  
1139 precise linguistic expression not to mention the assumption that learners have insight into  
1140 the structure of their own concepts. It is known that children's vocabularies differ from  
1141 adults', raising the concern that some of our results reflect language use rather than the  
1142 concepts being articulated. While our artificial environment contains only simple objects

and basic features that are familiar to even young children, there is evidence that children's speech does not distinguish as well among quantifier usage (e.g., all, each, every) until late in childhood (Brooks & Braine, 1996; Inhelder & Piaget, 1958). Thus, it could be that linguistic imprecision is behind some of the differences between children's and adults' guesses. For instance, this seems like a potential explanation for the lack of any exactly correct guesses from children about the quantifier-dependent rule 4 "exactly one is blue". However, a closer look at responses reveals that only 11/47 children guessed a rule that mentioned blue at all. Meanwhile 37/50 of adults' rules mentioned blue, but all but seven of these were wrong about the particulars of the quantification. In many cases other potential quantifications were not ruled out by adults' testing. For instance, several subjects never tried adding more than one blue object to a scene and later responded that *at least one* object must be blue. Thus, it seems that children's rules simply picked out different features of the scenes than adults. An interesting question is whether, in the cases where a child's guess is logically inconsistent with some of their learning data, this is because their representation itself is imprecise, or because their verbal description imprecisely describes their representation. Another possibility could be that adults are better introspectors than children, better able to "read out" the structure of their own representations (Morris, 2021). While these are intriguing possibilities our current experiment cannot fully resolve these explanations.

## Implications for theories of concepts

Psychological theories of concepts have oscillated between symbolic accounts—that seek to explain conceptual productivity and creativity—and similarity accounts—that seek to explain how concepts drive probabilistic generalization. The constructivist framework is based in the symbolic camp, however it inherits many of the advantages of similarity accounts by maintaining a relationship with probabilistic inference embodied by the stochastic mechanisms of generation and search. Thus, we see our findings as support for recent claims that higher level cognition utilizes some form of stochastic generative sampling to approximate rational inference (Bramley, Dayan, et al., 2017; Sanborn et al., 2021; Zhu, Sanborn, & Chater, 2020) and that this might also explain aspects of human cultural and technological development that take place over populations and multiple generations (Krafft, Shmueli, Griffiths, Tenenbaum, et al., 2021).

While neither the PCFG or IDG are oven-ready process models of human concept formation, they provide a useful starting point for thinking about process accounts. The PCFG framework describes normative inference in the limit of infinite sampling, but also provides a mechanism for both generating and adapting samples. The IDG is a hybrid that

1178 seeds hypotheses by trying to describe patterns that are present in observations rather  
 1179 than merely possible, making it more sample-efficient as a brute force approach to inference  
 1180 in situations where a learner already has some positive or demonstrative evidence of a  
 1181 concept. However its success is dependent on the learner generating or encountering scenes  
 1182 that exemplify and isolate causally relevant features. With enough evidence both  
 1183 approaches should favor the ground truth but with little evidence the PCFG will tend to  
 1184 entertain many concepts that the IDG does not.

1185 While the IDG captured the data better here, it is not a complete account because,  
 1186 even with instance-inspired starting point, we still need to explain how a learner adapts in  
 1187 light of new evidence. Following a number of recent research lines (Bramley, Mayrhofer,  
 1188 Gerstenberg, & Lagnado, 2017; Dasgupta, Schulz, & Gershman, 2017; Ullman, Goodman,  
 1189 & Tenenbaum, 2012), we see incremental mutation of one or a few focal hypotheses in the  
 1190 light of evidence as a promising approach. For instance, a learner might use an observation  
 1191 to generate an initial idea akin to our IDG, but then explore permutations to this to  
 1192 generate new scenes to test (Oaksford & Chater, 1994), and to account for these tests  
 1193 (Fränken et al., 2022). While older models like RULEX (Nosofsky & Palmeri, 1998;  
 1194 Nosofsky et al., 1994) provide candidate heuristics for achieving such a search over theories,  
 1195 their long run behavior lacks a clear relationship with computational-level rationality  
 1196 (Navarro, 2005). However, if a learners' adaptations approximate a valid approximation  
 1197 scheme, for instance accepting proposed permutations with the Metropolis-Hastings  
 1198 probability  $\max(1, \frac{P(h')}{P(h^t)})$  (Bramley, Dayan, et al., 2017; Dasgupta et al., 2016; Hastings,  
 1199 1970; Thaker et al., 2017), they can start to explain why more probable hypotheses are  
 1200 discovered more often as well as explaining probability matching and order effects are  
 1201 inevitable consequences of approximation (see Fränken et al., 2022). Since the endpoint of  
 1202 an MCMC search approaches an independent posterior sample, we would expect a  
 1203 population of such searchers to end up with a set of hypotheses that look like posterior  
 1204 samples. Moreover, since individual searchers have finite time to search, we would expect  
 1205 order effects and dependence in their ideas over time. To the extent that participants  
 1206 deviate from a probabilistically valid approximation scheme, for instance "hill climbing" or  
 1207 accepting only strictly better fitting ideas, we might also explain how they can get stuck in  
 1208 local optima and exhibit mal-adaptive order effects like garden paths (Gelpi, Prystawski,  
 1209 Lucas, & Buchsbaum, 2020). Taking the idea that earlier hypotheses carry information  
 1210 about older evidence and inference, we might also think of a population of such hypotheses  
 1211 as a kind of particle filter (Bramley, Dayan, et al., 2017; Daw & Courville, 2008). While  
 1212 acting primarily as a computational level norm, the PCFG prior provides useful  
 1213 infrastructure for hypothesis search. For example, prior production weights can be used to

1214 adapt an existing hypothesis by partially “regrowing” it (Goodman et al., 2008).  
1215 Furthermore, prior production weights implied by a generative prior mechanism combined  
1216 data likelihoods allows for the principled acceptance or rejection of new proposals in an  
1217 MCMC-like search scheme. For this to become a fully satisfying account of constructivist  
1218 inference this would need to be paired with a mechanism for scene generation in line with  
1219 those we sketch in Figure 3c&d, so explaining anchoring, order effects, probability  
1220 matching and confirmation bias in a unified account (Klahr & Dunbar, 1988).

1221 Our modeling of generalizations revealed that there is no straightforward family  
1222 resemblance between the features of rule-following training scenes (generated by the  
1223 participant) and rule-following generalization scenes (as pre-selected for the experiment).  
1224 This resulted in the Similarity model performing at chance and also being completely  
1225 uncorrelated with participants while all our symbolic model variants received support.  
1226 While this is far from an exhaustive comparison with sub-symbolic concept models, even a  
1227 successful similarity-driven account of generalizations would only account for half of the  
1228 behavior in this task. As well as generalizing systematically, participants gave detailed  
1229 natural language descriptions of their ideas. The majority of these we could convert into  
1230 logical statements (86%) that predicted most generalizations (72%: children, 84%: adults)  
1231 and were consistent with the majority of their learning data (71%: children, 87%: adults).  
1232 Any subsymbolic account of concepts would essentially need to be paired with an  
1233 explanation for *how* people generate these verbal descriptions of their non-symbolic  
1234 concepts that nonetheless reflect their use (cf. Dennett, 1988). Arguably, this task is no  
1235 easier than the one of generating a symbolic hypothesis about the nature of the world in  
1236 the first place. Thus we feel that our results are more straightforwardly explained by our  
1237 symbolic account whereby the logical structure of the hypotheses participants describe is  
1238 actually the causal mechanism driving their generalizations rather than some form of  
1239 computationally expensive but behaviorally impotent retrospective confabulation (cf.  
1240 Johansson et al., 2008). Our generalization analysis also showcases the difficulty of  
1241 predicting human behavior in a setting where there is such a large and long-tailed space of  
1242 similarly plausible rules an individual might be using to drive their generalizations.  
1243 Modeling symbolic inference directly from the learning input had some predictive power for  
1244 adults’ generalizations, but simply by asking participants for their best guess, we could  
1245 immediately get a far better handle on how they would generalize.

1246 While we did not provide a fully satisfying model of scene generation, we did show  
1247 that participant-generated scenes were better understood as adapting earlier scenes than as  
1248 being created from scratch. We argued that this is consistent with testing driven by one or  
1249 a couple of conceptually neighboring hypotheses, either generalizing their predictions or

1250 contrasting them. This is in some ways a return to pre-Bayesian ideas in philosophy of  
1251 science in testing permits falsification but not confirmation. Even when a hypothesis  $h$   
1252 survives repeated confirmatory tests, or repeated head-to-head challenges from local  
1253 alternatives, we might think of it as gaining a degree of confirmation, but there always  
1254 remains the specter of potential future falsification (cf. Popper, 1959). We think this better  
1255 reflects the state of a constructivist learner who cannot know, until discovering it, whether  
1256 some better hypothesis is waiting in the wings.

1257 For a learner limited to a few hypotheses at a time, the approach has clear virtues:  
1258 It links the process of adapting a hypotheses with that of coming up with new scenes to  
1259 test and links the outcome of tests to the subsequent inferential step of supplanting or  
1260 reinforcing the current favored hypothesis. Since learners are always reusing at least some  
1261 feature or other, it allows the learner's two tasks to support the other, with reuse of  
1262 modified previous tests and minimal positive examples minimizing the cognitive and  
1263 physical costs of generating both new tests and new hypotheses (Gershman & Niv, 2010).

## 1264 Implications for theories of development

1265 Our analyses revealed a variety of developmental differences. Children's guesses  
1266 were more complex than adults', and consequently we could capture them with a  
1267 significantly "flatter" generation process that inherently produced a wider diversity of  
1268 hypotheses. This is potentially normative: Having been exposed to less evidence, with less  
1269 idea what conceptual compositions and fragments will be useful in understanding their  
1270 environment, we should expect children's construction process to be less fine-tuned. In  
1271 other words, children are justified in entertaining a wider set of ideas than adults.  
1272 However, we noted there are several algorithmic stories that could underpin this diversity:  
1273 (1) children might simply have hypothesis generation mechanism that embodies a  
1274 rationally flatter latent prior, (2) they might additionally explore theory space more  
1275 radically, over and above differences in the relative credibility their latent prior actually  
1276 attaches to different possibilities (Gopnik, 2020; Lucas et al., 2014; Wu, Schulz,  
1277 Speekenbrink, Nelson, & Meder, 2018) or (3) we also considered that children's generation  
1278 mechanisms might be more dominated by "bottom-up" processes. We take our comparison  
1279 of PCFG and IDG to speak against option 3. Adults' hypotheses were, as far as we could  
1280 tell, at least as anchored to idiosyncratic patterns of their learning data as children's.  
1281 However, these data do not distinguish clearly between options (1) and (2). To do this, one  
1282 would need to measure children and adults' prior distributions directly. If children's  
1283 guesses shift within a problem in a way that is less sensitive to their own relative subjective  
1284 probabilities than adults, this would support the idea that children's hypothesis generation

is more “high temperature” exploratory than adults’ (Gopnik, 2020), over and above differences in the flatness of their latent prior. Importantly, while the endpoints of children’s theorizing were more diverse than adults’, the cognition required to produce their hypotheses still highly systematic. Children were able to implement a stable-enough symbolic generation or adaptation mechanism to produce meaningful symbolic hypotheses on the large majority of trials, referring to the features and relations they encountered. Even when their hypotheses did a poor job of explaining all the learning data, the hypothesis construction process did not break down entirely as it would if childlike brain activity were simply random and disorganized. However, the issue remains whether there is just more noise in children’s behavior—e.g., they are just a bit more easily distracted compared to adults—as opposed something like a greater inclination to explore.

Another aspect of constructivism that we did not focus on here but that is critical to understanding development, is the idea that over time, learners can chunk, cache and recursively reuse concepts to build ever richer ones (cf. Zhao, Bramley, & Lucas, 2022). As such the conceptual library of an adult ought to be more advanced, containing more powerful and complex concepts that can be readily reused to build new concepts. This might lead to a prediction of a different pattern of guesses than we found here. That is, we might have expected adults’ concepts to look more complex than children’s, not because they are built from more parts, but because the parts they are built from are, themselves, more complex. We suspect that the reason we did not find this sort of pattern here is that our task used very basic abstract features. Presumably our shape and geometric relation concepts are fairly established by around the age of 10. We predict that this would not hold in more applied domains where adults are able to draw on advanced concepts. For instance, when theorizing about economic conditions an adult might refer advanced primitives like “power laws”, “compound growth” or “arbitrage” that we would not expect to exist yet in the conceptual repertoire of many 9-11 year olds.

As well as producing more complex guesses, children also produced more elaborate scenes during learning. One possible characterization is that children’s active scene construction was more exploration-driven and less hypothesis-driven than adults’ (Wu et al., 2018), perhaps mixing more hypotheses-free exploration-driven actions in with hypothesis-driven systematic ones (Meder, Wu, Schulz, & Ruggeri, 2021). Indeed, differences in active exploration are the other side of the coin of the high temperature search idea (Friston et al., 2016; Gopnik, 2020; Klahr & Dunbar, 1988; E. Schulz, Klenske, Bramley, & Speekenbrink, 2017). However within each trial, children’s testing was more repetitive than adults’, suggesting that they made slower progress in exploring the problem space, or were generally less able to keep track of what they had done. The problem of

generating informative tests is not quite the same as that of finding the right hypothesis. It is important to avoid redundancy and, in combination, serve to test a wide variety of salient hypotheses. In this sense, adults' testing behavior was more systematic, better reducing global measures of uncertainty and potentially reflecting a more metacognitive control over learning (Kuhn & Brannock, 1977; Oaksford & Chater, 1994).

Curiously, children were more likely to refer to relational and positional properties in their guesses, while adults were most likely to make guesses that pertained to the primary object features (color and size). This is an independently interesting finding. Since relational features are structurally more complex than primitive features, we might have predicted they would be more readily evoked by adults. It could be that children bought in more to the scientific reasoning cover story, treating mechanistic explanations, such as that objects must touch or be positioned in particular ways to produce stars, as credible (Gelman, 2004). Conversely, adults may have been more likely to expect Gricean considerations to apply, e.g. that experimenters would likely set simple rules using salient but abstract features like color over perceptually ambiguous properties like position (Szollosi & Newell, 2020). However, it could also be the case that there are deeper differences between the experiences of children and adults that render structural features more relevant to children and surface features more relevant to adults.

Children's guesses were also less consistent with their evidence than adults'. This might be because they were less able to extract common features across all eight learning scenes (Ruggeri & Feufel, 2015; Ruggeri & Lombrozo, 2015). However, it could also be a consequence of a more generalized limitation in ability to generate, store and compare hypotheses. With a flatter prior and limited sampling, one has a lower chance of ever generating a hypothesis that can explain all the evidence. Children also under-generalized, often selecting only 1 or 2 of the 8 test scenes (there was actually always 4) doing so even when their symbolic guesses predicted more should be selected. It could be that children found this part of the task overwhelming, perhaps tending to stop after identifying one or two hypothesis consistent scenes rather than evaluating all of them. In sum, it seems children were less able to come up with a concise description of all the evidence generated, reflecting both a less developed metacognitive awareness and the skills needed (both verbal and conceptual) to extract patterns.

## Conclusions

We analyzed an experiment combining rich qualitative and quantitative measures of children's and adults' inductive inference. We found a number of developmental differences and demonstrated that we can make sense of these through a constructivist lens. Our

- <sub>1356</sub> results add empirical support and theoretical detail to recent characterizations of children  
<sub>1357</sub> as more diverse thinkers and active learners than adults, and bring us closer to a  
<sub>1358</sub> computational understanding of human learning across the lifespan.

1359

## References

- 1360 Allen, K. R., Smith, K. A., & Tenenbaum, J. B. (2020). Rapid trial-and-error learning  
1361 with simulation supports flexible tool use and physical reasoning. *Proceedings of the  
1362 National Academy of Sciences*, 117(47), 29302–29310.
- 1363 Bonawitz, E. B., Denison, S., Gopnik, A., & Griffiths, T. L. (2014). Win-stay, lose-sample:  
1364 A simple sequential algorithm for approximating Bayesian inference. *Cognitive  
1365 Psychology*, 74, 35–65.
- 1366 Box, G. E. (1976). Science and statistics. *Journal of the American Statistical Association*,  
1367 71(356), 791–799.
- 1368 Bramley, N. R., Dayan, P., Griffiths, T. L., & Lagnado, D. A. (2017). Formalizing  
1369 Neurath's ship: Approximate algorithms for online causal learning. *Psychological  
1370 Review*, 124(3), 301–338.
- 1371 Bramley, N. R., Jones, A., Gureckis, T. M., & Ruggeri, A. (2022). Changing many things  
1372 at once sometimes makes for a good experiment, and children know that.
- 1373 Bramley, N. R., Lagnado, D. A., & Speekenbrink, M. (2015). Conservative forgetful  
1374 scholars: How people learn causal structure through interventions. *Journal of  
1375 Experimental Psychology: Learning, Memory & Cognition*, 41(3), 708–731.
- 1376 Bramley, N. R., Mayrhofer, R., Gerstenberg, T., & Lagnado, D. A. (2017). Causal learning  
1377 from interventions and dynamics in continuous time. In *Proceedings of the 39<sup>th</sup>  
1378 Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive Science  
1379 Society.
- 1380 Bramley, N. R., Rothe, A., Tenenbaum, J. B., Xu, F., & Gureckis, T. M. (2018).  
1381 Grounding compositional hypothesis generation in specific instances. In *Proceedings  
1382 of the 40<sup>th</sup> Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive  
1383 Science Society.
- 1384 Brooks, P. J., & Braine, M. D. (1996). What do children know about the universal  
1385 quantifiers all and each? *Cognition*, 60(3), 235–268.
- 1386 Bruner, J. S., Goodnow, J. J., & Austin, G. A. (1956). *A study of thinking*. Routledge.
- 1387 Bruner, J. S., Jolly, A., & Sylva, K. (1976). *Play: Its role in development and evolution*.  
1388 Penguin.
- 1389 Campbell, D. T. (1960). Blind variation and selective retention in creative thought as in  
1390 other knowledge processes. *Psychological Review*, 67, 380–400.
- 1391 Carey, S. (1985). Are children fundamentally different kinds of thinkers and learners than  
1392 adults. *Thinking and Learning Skills*, 2, 485–517.
- 1393 Carey, S. (2009). *The origin of concepts: Oxford series in cognitive development*. Oxford  
1394 University Press, England.

- 1395 Chen, Z., & Klahr, D. (1999). All other things being equal: Acquisition and transfer of the  
1396 control of variables strategy. *Child Development*, 70(5), 1098–1120.
- 1397 Church, A. (1932). A set of postulates for the foundation of logic. *Annals of mathematics*,  
1398 346–366.
- 1399 Clark, A. (2012). Whatever next? predictive brains, situated agents, and the future of  
1400 cognitive science. *Behavioral Brain Sciences*, 1–86.
- 1401 Coenen, A., Rehder, R., & Gureckis, T. M. (2015). Strategies to intervene on causal  
1402 systems are adaptively selected. *Cognitive Psychology*, 79, 102–133.
- 1403 Dasgupta, I., Schulz, E., & Gershman, S. J. (2016). Where do hypotheses come from?  
1404 *Center for Brains, Minds and Machines (preprint)*.
- 1405 Dasgupta, I., Schulz, E., & Gershman, S. J. (2017). Where do hypotheses come from?  
1406 *Cognitive Psychology*, 96, 1–25.
- 1407 Daw, N., & Courville, A. (2008). The pigeon as particle filter. *Advances in neural  
1408 information processing systems*, 20, 369–376.
- 1409 Dennett, D. C. (1988). The intentional stance in theory and practice. In R. Byrne &  
1410 A. Whiten (Eds.), *Machiavellian intelligence* (pp. 180–202). Oxford, UK: Oxford  
1411 University Press.
- 1412 Dennett, D. C. (1991). *Consciousness explained*. London, UK: Penguin.
- 1413 Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., ... Tenenbaum, J. B.  
1414 (2020). Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep  
1415 bayesian program learning. *arXiv preprint arXiv:2006.08381*.
- 1416 Fedyk, M., & Xu, F. (2018). The epistemology of rational constructivism. *Review of  
1417 Philosophy and Psychology*, 9(2), 343–362.
- 1418 Feldman, J. (2000). Minimization of Boolean complexity in human concept learning.  
1419 *Nature*, 407(6804), 630.
- 1420 Fodor, J. A. (1975). *The language of thought* (Vol. 5). Harvard University Ppress.
- 1421 Fränken, J.-P., Theodoropoulos, N. C., & Bramley, N. R. (2022). Algorithms for  
1422 adaptation in inductive inference. *Cognitive Psychology*.
- 1423 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., Pezzulo, G., et al. (2016). Active  
1424 inference and learning. *Neuroscience & Biobehavioral Reviews*, 68, 862–879.
- 1425 Gelman, S. A. (2004). Psychological essentialism in children. *Trends in Cognitive Sciences*,  
1426 8(9), 404–409.
- 1427 Gelpi, R., Prystawski, B., Lucas, C. G., & Buchsbaum, D. (2020). Incremental hypothesis  
1428 revision in causal reasoning across development.
- 1429 Gershman, S. J., & Niv, Y. (2010). Learning latent structure: carving nature at its joints.  
1430 *Current Opinion in Neurobiology*, 20(2), 251–256.

- 1431 Gettys, C. F., & Fisher, S. D. (1979). Hypothesis plausibility and hypothesis generation.  
1432     *Organizational behavior and human performance*, 24(1), 93–110.
- 1433 Ginsburg, S. (1966). *The mathematical theory of context free languages*. McGraw-Hill  
1434     Book Company.
- 1435 Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational  
1436     analysis of rule-based concept learning. *Cognitive Science*, 32(1), 108–154.
- 1437 Goodman, N. D., Ullman, T. D., & Tenenbaum, J. B. (2011). Learning a theory of  
1438     causality. *Psychological Review*, 118(1), 110–9.
- 1439 Gopnik, A. (1996). The scientist as child. *Philosophy of Science*, 63(4), 485–514.
- 1440 Gopnik, A. (2020). Childhood as a solution to explore–exploit tensions. *Philosophical  
1441     Transactions of the Royal Society B*, 375(1803), 20190502.
- 1442 Gopnik, A., Glymour, C., Sobel, D., Schulz, L. E., Kushnir, T., & Danks, D. (2004). A  
1443     theory of causal learning in children: Causal maps and Bayes nets. *Psychological  
1444     Review*, 111, 1–31.
- 1445 Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources:  
1446     Levels of analysis between the computational and the algorithmic. *Topics in  
1447     Cognitive Science*, 7, 217–229.
- 1448 Griffiths, T. L., & Tenenbaum, J. B. (2009). Theory-based causal induction. *Psychological  
1449     Review*, 116, 661–716.
- 1450 Gureckis, T. M., & Markant, D. B. (2012, September). Self-Directed Learning: A  
1451     Cognitive and Computational Perspective. *Perspectives on Psychological Science*,  
1452     7(5), 464–481.
- 1453 Hastings, W. K. (1970). Monte carlo sampling methods using markov chains and their  
1454     applications.
- 1455 Heath, C. (2004). *Zendo—Design History*. Retrieved from  
1456     <http://www.koryheath.com/zendo/design-history/>
- 1457 Howson, C., & Urbach, P. (2006). *Scientific reasoning: the Bayesian approach*. Open  
1458     Court Publishing.
- 1459 Inhelder, B., & Piaget, J. (1958). *The growth of logical thinking from childhood to  
1460     adolescence: An essay on the construction of formal operational structures* (Vol. 22).  
1461     Psychology Press.
- 1462 Johansson, P., Hall, L., & Sikström, S. (2008). From change blindness to choice blindness.  
1463     *Psychologica*, 51(2), 142–155.
- 1464 Johnson, R. B., Onwuegbuzie, A. J., & Turner, L. A. (2007). Toward a definition of mixed  
1465     methods research. *Journal of Mixed Methods Research*, 1(2), 112–133.
- 1466 Johnson-Laird, P. N. (1983). *Mental models: Towards a cognitive science of language*,

- 1467         *inference, and consciousness*. Cambridge: Cambridge University Press.
- 1468     Kemp, C., & Tenenbaum, J. B. (2009). Structured statistical models of inductive  
1469             reasoning. *Psychological Review*, 116(1), 20.
- 1470     Klahr, D., & Dunbar, K. (1988). Dual space search during scientific reasoning. *Cognitive  
1471             Science*, 12(1), 1–48.
- 1472     Klahr, D., Fay, A. L., & Dunbar, K. (1993). Heuristics for scientific experimentation: A  
1473             developmental study. *Cognitive Psychology*, 25(1), 111–146.
- 1474     Klahr, D., Zimmerman, C., & Jirout, J. (2011). Educational interventions to advance  
1475             children's scientific thinking. *Science*, 333(6045), 971–975.
- 1476     Klayman, J., & Ha, Y.-w. (1989). Hypothesis testing in rule discovery: Strategy, structure,  
1477             and content. *Journal of Experimental Psychology: Learning, Memory & Cognition*,  
1478             15(4), 596.
- 1479     Komatsu, L. K. (1992). Recent views of conceptual structure. *Psychological bulletin*,  
1480             112(3), 500.
- 1481     Krafft, P. M., Shmueli, E., Griffiths, T. L., Tenenbaum, J. B., et al. (2021). Bayesian  
1482             collective learning emerges from heuristic social learning. *Cognition*, 212, 104469.
- 1483     Krippendorff, K. (2012). *Content analysis: An introduction to its methodology*. Sage.
- 1484     Kruschke, J. K. (1992). Alcove: an exemplar-based connectionist model of category  
1485             learning. *Psychological Review*, 99(1), 22.
- 1486     Kuhn, D., & Brannock, J. (1977). Development of the isolation of variables scheme in  
1487             experimental and “natural experiment” contexts. *Developmental Psychology*, 13(1),  
1488             9.
- 1489     Lagnado, D. A., & Sloman, S. A. (2006). Time as a guide to cause. *Journal of  
1490             Experimental Psychology: Learning, Memory & Cognition*, 32(3), 451–60.
- 1491     Lai, L., & Gershman, S. J. (2021). Policy compression: an information bottleneck in action  
1492             selection.
- 1493     Lakatos, I. (1976). Falsification and the methodology of scientific research programmes. In  
1494             *Can theories be refuted?* (pp. 205–259). Springer.
- 1495     Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building  
1496             machines that learn and think like people. *Behavioral and Brain Sciences*, 40.
- 1497     Lapidow, E., & Walker, C. M. (2020). The search for invariance: repeated positive testing  
1498             serves the goals of causal learning. *Language and concept acquisition from infancy  
1499             through childhood*, 197–219.
- 1500     Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and  
1501             reversals. In *Soviet physics doklady* (Vol. 10, pp. 707–710).
- 1502     Lewis, O., Perez, S., & Tenenbaum, J. (2014). Error-driven stochastic search for theories

- 1503 and concepts. In *Proceedings of the Annual Meeting of the Cognitive Science Society*  
1504 (Vol. 36).
- 1505 Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human  
1506 cognition as the optimal use of limited computational resources. *Behavioral and brain*  
1507 *sciences*, 43.
- 1508 Lieder, F., Griffiths, T. L., Huys, Q. J., & Goodman, N. D. (2018). The anchoring bias  
1509 reflects rational use of cognitive resources. *Psychonomic bulletin & review*, 25(1),  
1510 322–349.
- 1511 Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). Sustain: a network model of category  
1512 learning. *Psychological Review*, 111(2), 309.
- 1513 Lucas, C. G., Bridgers, S., Griffiths, T. L., & Gopnik, A. (2014). When children are better  
1514 (or at least more open-minded) learners than adults: Developmental differences in  
1515 learning the forms of causal relationships. *Cognition*, 131(2), 284–299.
- 1516 Lucas, C. G., & Griffiths, T. L. (2010). Learning the form of causal relationships using  
1517 hierarchical bayesian models. *Cognitive Science*, 34(1), 113–147.
- 1518 Luce, D. R. (1959). *Individual choice behavior*. New York: Wiley.
- 1519 Markant, D. B., & Gureckis, T. M. (2014). Is it better to select or to receive? learning via  
1520 active and passive hypothesis testing. *Journal of Experimental Psychology: General*,  
1521 143(1), 94.
- 1522 Marr, D. (1982). *Vision*. New York: Freeman & Co.
- 1523 McCormack, T., Bramley, N. R., Frosch, C., Patrick, F., & Lagnado, D. A. (2016).  
1524 Children's use of interventions to learn causal structure. *Journal of Experimental*  
1525 *Child Psychology*, 141, 1–22.
- 1526 Meder, B., Wu, C. M., Schulz, E., & Ruggeri, A. (2021). Development of directed and  
1527 random exploration in children. *Developmental Science*, 24(4), e13095.
- 1528 Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning.  
1529 *Psychological Review*, 85(3), 207.
- 1530 Meng, Y., Bramley, N., & Xu, F. (2018). Children's causal interventions combine  
1531 discrimination and confirmation. In *Proceedings of the 40<sup>th</sup> Annual Meeting of the*  
1532 *Cognitive Science Society*.
- 1533 Michalski, R. S. (1969). On the quasi-minimal solution of the general covering problem. In  
1534 *Proceedings of the 5<sup>th</sup> Annual Symposium on Information Processing* (Vol. A3, pp.  
1535 125–128).
- 1536 Morris, A. (2021). Invisible gorillas in the mind: Internal inattentional blindness and the  
1537 prospect of introspection training.
- 1538 Navarro, D. J. (2005). Analyzing the rulex model of category learning. *Journal of*

- 1539         *Mathematical Psychology*, 49(4), 259–275.
- 1540     Nelson, J. D., Divjak, B., Gudmundsdottir, G., Martignon, L. F., & Meder, B. (2014).  
1541             Children's sequential information search is sensitive to environmental probabilities.
- 1542             *Cognition*, 130(1), 74–80.
- 1543     Nickerson, R. S. (1998). Confirmation bias: A ubiquitous phenomenon in many guises.  
1544             *Review of General Psychology*, 2(2), 175.
- 1545     Nosofsky, R. M., & Palmeri, T. J. (1998). A rule-plus-exception model for classifying  
1546             objects in continuous-dimension spaces. *Psychonomic Bulletin & Review*, 5(3),  
1547             345–369.
- 1548     Nosofsky, R. M., Palmeri, T. J., & McKinley, S. C. (1994). Rule-plus-exception model of  
1549             classification learning. *Psychological Review*, 101(1), 53.
- 1550     Nye, M. I., Solar-Lezama, A., Tenenbaum, J. B., & Lake, B. M. (2020). Learning  
1551             compositional rules via neural program synthesis. *arXiv preprint arXiv:2003.05562*.
- 1552     Oaksford, M., & Chater, N. (1994). Another look at eliminative and enumerative behaviour  
1553             in a conceptual task. *European Journal of Cognitive Psychology*, 6(2), 149–169.
- 1554     Oaksford, M., & Chater, N. (2007). *Bayesian rationality: the probabilistic approach to  
1555             human reasoning*. Oxford: Oxford University Press.
- 1556     Osborne, M., Garnett, R., Ghahramani, Z., Duvenaud, D. K., Roberts, S. J., & Rasmussen,  
1557             C. (2012). Active learning of model evidence using bayesian quadrature. *Advances in  
1558             neural information processing systems*, 25.
- 1559     Phillips, D. C. (1995). The good, the bad, and the ugly: The many faces of constructivism.  
1560             *Educational researcher*, 24(7), 5–12.
- 1561     Piaget, J. (2013). *The construction of reality in the child* (Vol. 82). Routledge.
- 1562     Piaget, J., & Valsiner, J. (1930). *The child's conception of physical causality*. Transaction  
1563             Pub.
- 1564     Piantadosi, S. T. (2021). The computational origin of representation. *Minds and  
1565             Machines*, 31(1), 1–58.
- 1566     Piantadosi, S. T., & Jacobs, R. A. (2016). Four problems solved by the probabilistic  
1567             language of thought. *Current Directions in Psychological Science*, 25(1), 54–59.
- 1568     Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2012). Bootstrapping in a  
1569             language of thought: A formal model of numerical concept learning. *Cognition*,  
1570             123(2), 199–217.
- 1571     Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2016). The logical primitives of  
1572             thought: Empirical foundations for compositional cognitive models. *Psychological  
1573             Review*, 123(4), 392.
- 1574     Popper, K. (1959). *The logic of scientific discovery*. Routledge.

- 1575 Posner, M. I., & Keele, S. W. (1968). On the genesis of abstract ideas. *Journal of*  
1576 *Experimental Psychology*, 77(3p1), 353.
- 1577 Quine, W. v. O. (1969). *Word and object*. MIT press.
- 1578 Rothe, A., Lake, B. M., & Gureckis, T. M. (2017). Question asking as program generation.  
1579 In *Neural Information Processing Systems*.
- 1580 Ruggeri, A., & Feufel, M. (2015). How basic-level objects facilitate question-asking in a  
1581 categorization task. *Frontiers in Psychology*, 6, 918.
- 1582 Ruggeri, A., & Lombrozo, T. (2014). Learning by asking: How children ask questions to  
1583 achieve efficient search. In *Proceedings of the 36<sup>th</sup> annual meeting of the cognitive*  
1584 *science society* (pp. 1335–1340). Austin, TX: Cognitive Science Society.
- 1585 Ruggeri, A., & Lombrozo, T. (2015). Children adapt their questions to achieve efficient  
1586 search. *Cognition*, 143, 203–216.
- 1587 Ruggeri, A., Lombrozo, T., Griffiths, T. L., & Xu, F. (2016). Sources of developmental  
1588 change in the efficiency of information search. *Developmental Psychology*, 52(12),  
1589 2159.
- 1590 Ruis, L., Andreas, J., Baroni, M., Bouchacourt, D., & Lake, B. M. (2020). A benchmark  
1591 for systematic generalization in grounded language understanding. *arXiv preprint*  
1592 *arXiv:2003.05161*.
- 1593 Rule, J. S., Schulz, E., Piantadosi, S. T., & Tenenbaum, J. B. (2018). Learning list  
1594 concepts through program induction. *BioRxiv*, 321505.
- 1595 Rule, J. S., Tenenbaum, J. B., & Piantadosi, S. T. (2020). The child as hacker. *Trends in*  
1596 *Cognitive Sciences*.
- 1597 Sanborn, A. N., & Chater, N. (2016). Bayesian brains without probabilities. *Trends in*  
1598 *Cognitive Sciences*.
- 1599 Sanborn, A. N., Zhu, J., Spicer, J., Sundh, J., León-Villagrá, P., & Chater, N. (2021).  
1600 Sampling as the human approximation to probabilistic inference.
- 1601 Schulz, E., Klenske, E. D., Bramley, N. R., & Speekenbrink, M. (2017). Strategic  
1602 exploration in human adaptive control. In *Proceedings of the 39<sup>th</sup> Annual Meeting of*  
1603 *the Cognitive Science Society*. The Cognitive Sience Society.
- 1604 Schulz, L. E., Goodman, N. D., Tenenbaum, J. B., & Jenkins, A. C. (2008). Going beyond  
1605 the evidence: Abstract laws and preschoolers' responses to anomalous data.  
1606 *Cognition*, 109(2), 211–223.
- 1607 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2),  
1608 461–464.
- 1609 Shackle, S. (2015). Science and serendipity: famous accidental discoveries: Most scientific  
1610 breakthroughs take years of research—but often, serendipity provides the final push,

- 1611 as these historic discoveries show. *New Humanist*, 2.
- 1612 Shanks, D. R., Tunney, R. J., & McCarthy, J. D. (2002). A re-examination of probability  
1613 matching and rational choice. *Journal of Behavioral Decision Making*, 15(3),  
1614 233–250.
- 1615 Shepard, R. N. (1987). Toward a universal law of generalization for psychological science.  
1616 *Science*, 237(4820), 1317–1323.
- 1617 Shepard, R. N., & Chang, J.-J. (1963). Stimulus generalization in the learning of  
1618 classifications. *Journal of Experimental Psychology*, 65(1), 94.
- 1619 Sim, Z. L., & Xu, F. (2017). Learning higher-order generalizations through free play:  
1620 Evidence from 2-and 3-year-old children. *Developmental Psychology*, 53(4), 642.
- 1621 Simon, H. A. (2013). *Administrative behavior*. Simon and Schuster.
- 1622 Sobel, D. M., & Kushnir, T. (2006). The importance of decision making in causal learning  
1623 from interventions. *Memory & Cognition*, 34(2), 411–419.
- 1624 Stewart, N., Chater, N., & Brown, G. D. A. (2006). Decision by sampling. *Cognitive  
1625 Psychology*, 53(1), 1–26.
- 1626 Steyvers, M., Tenenbaum, J. B., Wagenmakers, E., & Blum, B. (2003). Inferring causal  
1627 networks from observations and interventions. *Cognitive Science*, 27, 453–489.
- 1628 Szollosi, A., & Newell, B. R. (2020). People as intuitive scientists: Reconsidering statistical  
1629 explanations of decision making. *Trends in Cognitive Sciences*.
- 1630 Tenenbaum, J. B. (1999). *A Bayesian framework for concept learning* (Unpublished  
1631 doctoral dissertation). Massachusetts Institute of Technology.
- 1632 Thaker, P., Tenenbaum, J. B., & Gershman, S. J. (2017). Online learning of symbolic  
1633 concepts. *Journal of Mathematical Psychology*, 77, 10–20.
- 1634 Turing, A. M. (1937). On computable numbers, with an application to the  
1635 entscheidungsproblem. *Proceedings of the London mathematical society*, 2(1),  
1636 230–265.
- 1637 Turing, A. M. (2009). Computing machinery and intelligence. In *Parsing the turing test*  
1638 (pp. 23–65). Springer.
- 1639 Tversky, A. (1977). Features of similarity. *Psychological Review*, 84(4), 327.
- 1640 Ullman, T. D., Goodman, N. D., & Tenenbaum, J. B. (2012). Theory learning as stochastic  
1641 search in the language of thought. *Cognitive Development*, 27(4), 455–480.
- 1642 Valkov, L., Chaudhari, D., Srivastava, A., Sutton, C., & Chaudhuri, S. (2018). Houdini:  
1643 Lifelong learning as program synthesis. In *Advances in Neural Information  
1644 Processing Systems* (pp. 8687–8698).
- 1645 Van Laarhoven, P. J., & Aarts, E. H. (1987). Simulated annealing. In *Simulated annealing:  
1646 Theory and applications* (pp. 7–15). Springer.

- 1647 Van Rooij, I., Blokpoel, M., Kwisthout, J., & Wareham, T. (2019). *Cognition and*  
1648        *intractability: A guide to classical and parameterized complexity analysis*. Cambridge  
1649        University Press.
- 1650 von Humboldt, W. (1863/1988). *On language*. New York: Cambridge University Press.
- 1651 Vul, E., Goodman, N. D., Griffiths, T. L., & Tenenbaum, J. B. (2009). One and done?  
1652        optimal decisions from very few samples. In *Proceedings of the 31<sup>st</sup> Annual Meeting*  
1653        *of the Cognitive Science Society* (Vol. 1, pp. 66–72). Austin, TX: Cognitive Science  
1654        Society.
- 1655 Wason, P. C. (1960). On the failure to eliminate hypotheses in a conceptual task.  
1656        *Quarterly Journal of Experimental Psychology*, 12(3), 129–140.
- 1657 Wason, P. C. (1968). Reasoning about a rule. *The Quarterly Journal of Experimental*  
1658        *Psychology*, 20(3), 273–281.
- 1659 Wu, C. M., Schulz, E., Speekenbrink, M., Nelson, J. D., & Meder, B. (2018).  
1660        Generalization guides human exploration in vast decision spaces. *Nature Human*  
1661        *Behaviour*, 2(12), 915–924.
- 1662 Xu, F. (2019). Towards a rational constructivist theory of cognitive development.  
1663        *Psychological Review*, 126(6), 841.
- 1664 Zhao, B., Bramley, N. R., & Lucas, C. (2022). Powering up causal generalization: A model  
1665        of human conceptual bootstrapping with adaptor grammars. In *Proceedings of the*  
1666        *44<sup>th</sup> annual meeting of the cognitive science society* (Vol. 44).
- 1667 Zhao, B., Lucas, C. G., & Bramley, N. R. (2022). How do people generalize causal  
1668        relations over objects? a non-parametric bayesian account. *Computational Brain &*  
1669        *Behavior*, 5(1), 22–44.
- 1670 Zhu, J.-Q., Sanborn, A. N., & Chater, N. (2020). The Bayesian sampler: Generic Bayesian  
1671        inference causes incoherence in human probability judgments. *Psychological Review*,  
1672        127(5), 719.

1673

## Appendix A: Models

### 1674 Generating PCFG model predictions

1675 We created a grammar (specifically a *probabilistic context free grammar* or PCFG;  
 1676 Ginsburg, 1966) that can be used to produce any rule that can be expressed with  
 1677 first-order logic and lambda abstraction referring to the features participants referred to in  
 1678 our task. The grammatical primitives we assumed are detailed in 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...   | $N_{\{<, >, =\}}(\lambda x_i : \mathcal{X}, N)$  |
| 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).

1679

There are multiple ways to implement a PCFG. Here we adopt a common approach

1680 to set up a set of string-rewrite rules (Goodman et al., 2008). Thus, each hypothesis begins  
 1681 life as a string containing a single *non-terminal symbol* (here,  $S$ ) that is replaced using

1682 rewrite rules, or *productions*. These productions are repeatedly applied to the string,  
 1683 replacing non-terminal symbols with a mixture of other non-terminal symbols and terminal  
 1684 fragments of first order logic, until no non-terminal symbols remain. The productions are  
 1685 so designed that the resulting string is guaranteed to be a valid grammatical expression  
 1686 and all grammatical expressions have a nonzero chance of being produced. In addition, by  
 1687 having the productions tie the expression to bound variables and truth statements, our  
 1688 PCFG serves as an automatic concept generator. Table A-2 details the PCFG we used in  
 1689 the paper.

1690 We use capital letters as non-terminal symbols and each rewrite is sampled from the  
 1691 available productions for a given symbol.<sup>11</sup> Because some of the productions involve  
 1692 branching (e.g.,  $B \rightarrow H(B, B)$ ), the resultant string can become arbitrarily long and  
 1693 complex, involving multiple boolean functions and complex relationships between bound  
 1694 variables.

1695 We include a variant that samples uniformly from the set of possible replacements  
 1696 in each case, but we also reverse engineer a set of productions that produce exactly the  
 1697 statistics of the empirical samples, as described in the main text.

1698 We used the process described in A-2 to produce a sample of 10,000 with a uniform  
 1699 generation prior and an additional 10,000 for each participant with a “held out”  
 1700 age-consistent prior based on the rule guesses of other participants in the requisite  
 1701 agegroup. For the flipped prior analyses, we used the sample generated for the  
 1702 chronologically first participant from the other agegroup.

### 1703 Generating instance driven (IDG) model predictions

1704 We used the algorithm proposed in Bramley et al. (2018) to produce a sample of  
 1705 10,000 “grounded hypotheses” for each participant and trial, splitting these evenly across  
 1706 the 8 learning scenes that participant produced and tested. For each, we generated two  
 1707 sets: One using a uniform construction weights, and one with an age-appropriate “held  
 1708 out” set of weights based on the rule guesses of other participants in the requisite agegroup.  
 1709 For the flipped prior analyses, we used the weights from the chronologically first participant  
 1710 from the other agegroup to generate samples inspired by the current participants’ evidence.

1711 To generate hypotheses as candidates for the hidden rule, the model uses the  
 1712 following procedure with probabilities either set to uniform or drawn from the PCFG-fitted

---

<sup>11</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. Rothe, Lake, & Gureckis, 2017).

**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                        | $C \rightarrow$  | $= (x_i, D1)$                          | $I(x_i, D2)$                           | $= (x_i, x_j, E1)^a$                  |
|                                 |                  | $I(x_i, x_j, E2)^a$                    | $\Gamma(x_i, x_j, E3)^a$               |                                       |
| Feature/value<br>(numeric only) | $D1 \rightarrow$ | value,                                 | feature                                |                                       |
|                                 | $D2 \rightarrow$ | value,                                 | feature                                |                                       |
| Feature<br>(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.

<sup>1713</sup> productions for adults or for children (Figure 7) and denoted with square brackets:

<sup>1714</sup> 1. **Observe.** either:

- <sup>1715</sup> (a) With probability  $[A \rightarrow B]$ : Sample a cone from the observation, then sample  
<sup>1716</sup> one of its features  $f$  with probability  $[G \rightarrow f]$ —e.g.,  $\{\#1\}$ :<sup>12</sup> “medium, size” or  
<sup>1717</sup>  $\{\#3\}$ : “red, color”.
- <sup>1718</sup> (b) With probability  $[A \rightarrow \text{Start}]$ : Sample two cones uniformly without replacement  
<sup>1719</sup> from the observation, and sample any shared or pairwise feature—e.g.,  
<sup>1720</sup>  $\{\#1, \#2\}$ : “size”, or “contact”

<sup>1721</sup> 2. **Functionize.** Bind a variable for each sampled cone in Step 1 and sample a true  
<sup>1722</sup> (in)equality statement relating the variable(s) and feature:

- <sup>1723</sup> (a) For a statement involving an unordered feature there is only one  
<sup>1724</sup> possibility—e.g.,  $\{\#3\}$ : “ $= (x_1, \text{red}, \text{color})$ ”, or for  $\{\#1, \#2\}$ : “ $= (x_1, x_2, \text{color})$ ”
- <sup>1725</sup> (b) For a single cone and an ordered feature, this could also be a nonstrict  
<sup>1726</sup> inequality ( $\geq$  or  $\leq$ ). We assume a learner only samples an inequality if it  
<sup>1727</sup> expands the number of cones picked out from the scene relative to an

<sup>12</sup> Numbers prepended with # refer to the labels on the cones in the example observation in Figure 2b.

1728 equality—e.g., in Figure 2b in the main text, there is also a large cone  $\{\#1\}$  so  
 1729 either  $\geq(x_1, \text{medium}, \text{size})$  or  $=x_1, \text{medium}, \text{size}$ ) might be selected with  
 1730 uniform probability.

- 1731 (c) For two cones and an ordered feature, either strict or non-strict inequalities  
 1732 could be sampled if the cones differ on the sampled feature, equivalently either  
 1733 equality or non-strict inequality could be selected if the cones do not differ on  
 1734 that dimension—e.g.,  $\{\#1,\#2\} > (x_1, x_2, \text{size})$ , or  $\{\#3,\#4\} \geq (x_1, x_2, \text{size})$ . In  
 1735 each case, the production weights from Figure 7 for the relevant completions are  
 1736 normalized and used to select the option.

- 1737 3. **Extend.** With probability  $\frac{[B \rightarrow D]}{[B \rightarrow D] + [B \rightarrow C(B, B)]}$  go to Step 4, otherwise sample a  
 1738 conjunction with probability  $[C(B, B) \rightarrow \text{And}]$  or a disjunction with probability  
 1739  $[C(B, B) \rightarrow \text{Or}]$  and repeat. For statements with two bound variables, Step 3 is  
 1740 performed for  $x_1$ , then again for  $x_2$ :

- 1741 (a) **Conjunction.** A cone is sampled from the subset picked out by the statement  
 1742 thus far and one of its features sampled with probability  $[G \rightarrow f]$ —e.g.,  $\{\#1\}$   
 1743  $\wedge (= (x_1, \text{green}, \text{color}), \geq (x_1, \text{medium}, \text{size}))$ . Again, inequalities are sampleable  
 1744 only if they increase the true set size relative to equality—e.g.,  
 1745 “ $\wedge (\leq (x_1, 3, \text{xposition}), \geq (x_1, \text{medium}, \text{size}))$ ”, which picks out more objects  
 1746 than “ $\wedge (= (x_1, 3, \text{xposition}), \geq (x_1, \text{medium}, \text{size}))$ ”.

- 1747 (b) **Disjunction.** An additional feature-value pair is selected uniformly from *either*  
 1748 unselected values of the current feature, *or* from a different feature—e.g.,  
 1749  $\vee (= (x_1, \text{color}, \text{red}), = (x_1, \text{color}, \text{blue}))$  or  $\vee (= (x_1, \text{color}, \text{blue}), \geq (x_1, \text{size}, 2))$ .  
 1750 This step is skipped if the statement is already true of all the cones in the  
 1751 scene.<sup>13</sup>

- 1752 4. **Flip.** If the inspiration scene is not rule following wrap the expression in a  $\neg()$ .

- 1753 5. **Quantify.** Given the contained statement, select true quantifier(s):

- 1754 (a) For statements involving a single bound variable (i.e., those inspired by a single  
 1755 cone in Step 1) the possible quantifiers simply depend on the number of the  
 1756 cones in the scene for which the statement holds. If the statement is true of all  
 1757 cones in the scene Quantifier is selected using probabilities  $[\text{Start} \rightarrow]$  combined  
 1758 with  $[L \rightarrow]$  where appropriate. If it is true of only a subset of the cones then

---

<sup>13</sup> We rounded positional features to one decimal place in evaluating rules to allow for perceptual uncertainty.

1759        $\forall(\lambda x_i : A, \mathcal{X})$  is censored and the probabilities re-normalized.  $K$  is set to match  
 1760       number of cones for which the statement is true.

- 1761       (b) Statements involving two bound variables in lambda calculus have two nested  
 1762       quantifier statements each selected as in (a). The inner statement quantifying  $x_2$   
 1763       is selected first based on truth value of the expression while taking  $x_1$  to refer to  
 1764       the cone observed in ‘1.’. The truth of the selected inner quantified statement is  
 1765       then assessed for all cones to select the outer quantifier—e.g.,  $\{\#3, \#4\}$   
 1766       “ $\wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size}))$ ” might become  
 1767       “ $\forall(\lambda x_1 : \exists(\lambda x_2 : \wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size})), \mathcal{X}), \mathcal{X})$ ”. The inner  
 1768       quantifier  $\exists$  is selected (three of the four cones are green  $\{\#1, \#2, \#4\}$ ), and  
 1769       the outer quantifier  $\forall$  is selected (all cones are less than or equal in size to a  
 1770       green cone).

1771       Note that a procedure like the one laid out above is, in principle, capable of  
 1772       generating any rule generated by the PCFG in Figure 7a&7b, but will only do so when  
 1773       exposed to an observation that exemplifies that rule, and will do so more often when the  
 1774       observation is inconsistent with as many other rules as possible (i.e., a minimal positive  
 1775       example). Step 4. allows that non-rule following scenes can be used to inspire rules  
 1776       involving a negation, for instance that “something is not upright” – which is semantically  
 1777       equivalent to saying that “nothing is upright”. Basing hypotheses on instances may  
 1778       improve the quality of the effective sample of hypotheses that the learner generates.

1779       One way to think of the IDG procedure is as a partial inversion of a PCFG. As  
 1780       illustrated by the blue text in the examples in Figure 2b in the main text. While the  
 1781       PCFG starts at the outside and works inward, the IDG starts from the central content and  
 1782       works outward out to a quantified statement, ensuring at each step that this final  
 1783       statement is true of the scene.

1784       We note that it is possible, in principle, to calculate a lower bound on the prior  
 1785       probability for the PCFG or IDG generating a hypothesis that a participant reported, even  
 1786       if it does not occur in our sample. This can be achieved by reverse engineering the  
 1787       production steps that would be needed to produce the precise encoded syntax. This is a  
 1788       lower bound because it does not count semantically equivalent “phrasings” of the  
 1789       hypothesis that e.g. mention features in different orders or use logically equivalent  
 1790       combinations of booleans. We found that complex expressions tend to have a large number  
 1791       of “phrasings”. In our sample-based approximation we implicitly treat semantically  
 1792       equivalent expressions as constituting the same hypothesis but note that determining  
 1793       semantic equivalence is an nontrivial aspect of constructivist inference that we do not fully

1794 address here.

1795 **Reverse engineering production child-like and adult-like production weights**

1796 To roughly accommodate the fact that each guess is based on different learning  
 1797 data, we regularized these counts by including a prior pseudo-count of 5 on all productions.  
 1798 This value was not fit to the data, and simply serves to smooth the predictions a little. For  
 1799 example, children's rules involved  $\exists$  263 times,  $\forall$  108 times and  $N$  297 times, so we  
 1800 assumed prior production weights of  
 1801  $\{263 + 5, 108 + 5, 297 + 5\}/(263 + 108 + 297 + 15) = \{.39, .17, .44\}$ . To avoid double  
 1802 counting the data in modeling subjects' specific guesses, we created a separate  
 1803 agegroup-appropriate prior production weighting for each participant based on the guesses  
 1804 of the other participants' from the same agegroup, but omitting their own guesses.

1805 **Appendix B - Model fitting details**

1806 **Full generalization model fits**

1807 As described in main text, we fit 18 model variants to participant's data. All models  
 1808 have between 0 and 2 parameters. For each model, we fit the parameter(s) by maximizing  
 1809 the model's likelihood of producing the participant data, using R's `optim` function. We  
 1810 compare models using the Bayesian Information Criterion (Schwarz, 1978) to accommodate  
 1811 their different numbers of fitted parameters.<sup>14</sup> Full results are in Table A-3.

1812 **Scene generation model fits**

1813 We used a grid search in increments of 0.05 to optimize  $\eta$  and  $\theta$  and directly  
 1814 optimized  $\lambda$  for each setting of  $\eta$  and  $\theta$ .

1815 **Appendix B: Free response coding**

1816 To analyze the free responses, we first had two coders go through all responses and  
 1817 categorize them as either:

---

<sup>14</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 Agegroup 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 since the held out priors are fit to different data than that being modeled. 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. Bias                          | children        | -1218.96        | 2445.47        | 0.32        |             | <b>16</b> | 50%      |
| 3. PCFG Uniform                  | children        | -1319.72        | 2647.00        |             | 58.17       | 0         | 61%      |
| 4. PCFG Uniform + Bias           | children        | -1208.93        | 2432.97        | 0.35        | 2.18        | 0         |          |
| 5. PCFG Flipped                  | children        | -1318.46        | 2644.47        |             | 8.97        | 1         | 66%      |
| 6. PCFG Flipped + Bias           | children        | -1207.28        | 2429.67        | 0.34        | 2.07        | 0         |          |
| 7. PCFG Agegroup                 | children        | -1319.58        | 2646.71        |             | 24.17       | 1         | 63%      |
| 8. PCFG Agegroup + Bias          | children        | -1208.63        | 2432.36        | 0.35        | 2.15        | 0         |          |
| 9. IDG Uniform                   | children        | -1298.73        | 2605.02        |             | 1.78        | 1         | 65%      |
| 10. IDG Uniform + Bias           | children        | -1193.90        | 2402.90        | 0.32        | 1.19        | 0         |          |
| 11. IDG Flipped                  | children        | -1315.49        | 2638.54        |             | 4.35        | 1         | 66%      |
| 12. IDG Flipped + Bias           | children        | -1199.22        | 2413.54        | 0.35        | 1.38        | 0         |          |
| 13. IDG Agegroup                 | children        | -1308.05        | 2623.65        |             | 2.51        | 2         | 69%      |
| 14. IDG Agegroup + Bias          | children        | -1193.41        | 2401.93        | 0.34        | 1.19        | 0         |          |
| 15. Similarity                   | children        | -1316.44        | 2640.42        |             | -1.99       | 0         | 41%      |
| 16. Similarity + Bias            | children        | -1214.71        | 2444.52        | 0.32        | -1.30       | 1         |          |
| 17. Symbolic Guess               | children        | -1143.69        | 2294.92        |             | 1.02        | 15        | 62%      |
| <b>18. Symbolic Guess + Bias</b> | <b>children</b> | <b>-1067.18</b> | <b>2149.47</b> | <b>0.26</b> | <b>0.80</b> | 9         |          |
| 1. Baseline                      | adults          | -1386.29        | 2772.59        |             |             | 2         | 50%      |
| 2. Bias                          | adults          | -1364.90        | 2737.40        | 0.15        |             | 6         | 50%      |
| 3. PCFG Uniform                  | adults          | -1320.64        | 2648.89        |             | 1.27        | 0         | 63%      |
| 4. PCFG Uniform + Bias           | adults          | -1253.52        | 2522.25        | 0.26        | 0.68        | 0         |          |
| 5. PCFG Flipped                  | adults          | -1294.91        | 2597.42        |             | 1.06        | 1         | 66%      |
| 6. PCFG Flipped + Bias           | adults          | -1229.18        | 2473.55        | 0.24        | 0.63        | 0         |          |
| 7. PCFG Agegroup                 | adults          | -1266.96        | 2541.51        |             | 0.94        | 1         | 69%      |
| 8. PCFG Agegroup + Bias          | adults          | -1203.64        | 2422.47        | 0.23        | 0.59        | 0         |          |
| 9. IDG Uniform                   | adults          | -1228.21        | 2464.02        |             | 0.67        | 2         | 69%      |
| 10. IDG Uniform + Bias           | adults          | -1179.12        | 2373.44        | 0.20        | 0.48        | 0         |          |
| 11. IDG Flipped                  | adults          | -1245.56        | 2498.72        |             | 0.76        | 0         | 73%      |
| 12. IDG Flipped + Bias           | adults          | -1179.23        | 2373.65        | 0.24        | 0.48        | 0         |          |
| 13. IDG Agegroup                 | adults          | -1188.28        | 2384.17        |             | 0.62        | 2         | 74%      |
| 14. IDG Agegroup + Bias          | adults          | -1134.58        | 2284.37        | 0.20        | 0.44        | 0         |          |
| 15. Similarity                   | adults          | -1359.05        | 2725.70        |             | -0.73       | 0         | 37%      |
| 16. Similarity + Bias            | adults          | -1337.55        | 2690.30        | 0.14        | -0.61       | 0         |          |
| 17. Symbolic Guess               | adults          | -893.49         | 1794.58        |             | 0.56        | <b>32</b> | 70%      |
| <b>18. Symbolic Guess + Bias</b> | <b>adults</b>   | <b>-880.59</b>  | <b>1776.38</b> | <b>0.08</b> | <b>0.50</b> | 4         |          |

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

- 1818 1. Correct: The subject gives exactly the correct rule or something logically equivalent
- 1819 2. Overcomplicated: The subject gives a rule that over-specifies the criteria needed to produce stars relative to the ground truth. This means the rule they give is logically sufficient but not necessary. For example, stipulating that "there must be a small red" is overcomplicated if the true rule is "there must be a red" because a scene could contain a medium or large red and emit stars.
- 1820 3. Overliberal: The opposite of overcomplicated. The subject gives a rule that under-specifies what must happen for the scene to produce stars. For example,
- 1821 4. Stipulative: The subject gives a rule that is logically equivalent to the true rule but adds unnecessary constraints. For example, stipulating that "there must be a blue" if the true rule is that "exactly one is blue".
- 1822 5. Ambiguous: The subject gives a rule that is logically equivalent to the true rule but is ambiguous about what it means. For example, stipulating that "there must be a red" when the true rule is "there must be a red" but the subject is not clear about what "red" means.
- 1823 6. Inconsistent: The subject gives a rule that contradicts the true rule. For example, stipulating that "there must be a red" when the true rule is "there must be a red" but the subject also stipulates that "there must be a blue".
- 1824 7. Uninformative: The subject gives a rule that does not provide any useful information about the scene. For example, stipulating that "there must be a red" when the true rule is "there must be a red" but the subject also stipulates that "there must be a blue".
- 1825 8. Unintelligible: The subject gives a rule that is not understandable or makes no sense. For example, stipulating that "there must be a red" when the true rule is "there must be a red" but the subject also stipulates that "there must be a blue".
- 1826 9. Unpredictable: The subject gives a rule that is not predictable or follows a pattern. For example, stipulating that "there must be a red" when the true rule is "there must be a red" but the subject also stipulates that "there must be a blue".

1827      This is logically necessary but not sufficient because a scene could contain blue  
 1828      objects but not produce stars because there is not exactly one of them.

- 1829      4. Different: The subject gives a rule that is intelligible but different from the ground  
 1830      truth in that it is neither necessary or sufficient for determining whether a scene will  
 1831      produce stars.
- 1832      5. Vague or multiple. Nuisance category.
- 1833      6. No rule. The subject says they cannot think of a rule.

1834      We were able to encode 205/238 (86%) of the children's responses and (219/250)  
 1835      87% for adults as correct, overcomplicated, overliberal or different. Table A-4 shows the  
 1836      complete confusion matrix. The two coders agreed 85% of the time, resulting in a Cohen's  
 1837      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> |

1838      We then had one coder familiar with the grammar go through each free response  
 1839      that was not assigned vague or no rule, and encode it as a function in our grammar. The  
 1840      second coder then blind spot checked 15% of these rules (64) and agreed in 95% of cases  
 1841      61/64. The 6 cases of disagreement were discussed and resolved. In 5/6 cases, this was in  
 1842      favor of the primary coder. The full set of free text responses along with the requisite  
 1843      classification, encoded rules are available in the [Online Repository](#).

#### 1844      Appendix C: Scene similarity measurement

1845      To establish the overall similarity between two scenes, we need to map the objects  
 1846      in a given scene to the objects in another scene (for example between the scenes in  
 1847      FigureA-1 a and b) and establish a reasonable cost for the differences between objects  
 1848      across dimensions. We also need a procedure for cases where there are objects in one scene

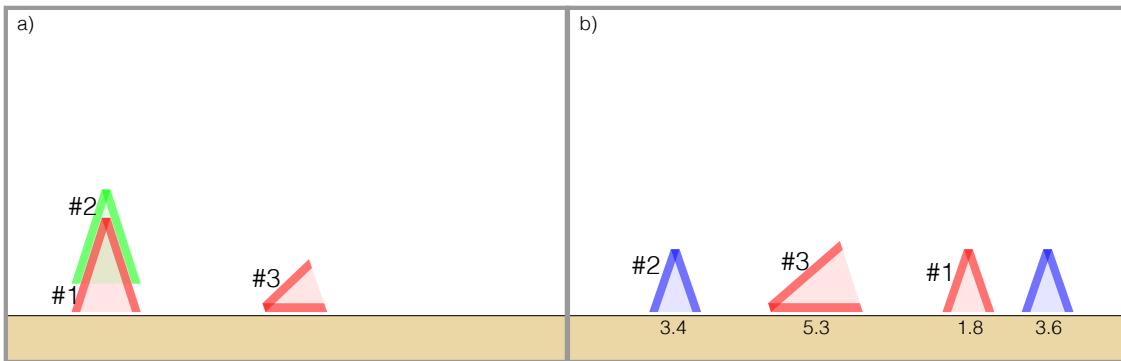
1849 that have no analogue in the other. We approach the calculation of similarity via the  
 1850 principle of minimum edit distance (Levenshtein, 1966). This means summing up the  
 1851 elementary operations required to convert scene (a) into scene (b) or visa versa. We assume  
 1852 objects can be adjusted in one dimension at a time (i.e. moving them on the  $x$  axis,  
 1853 rotating them, or changing their color, and so on.

1854 Before focusing on how to map the objects between the scenes we must decide how  
 1855 to measure the adjustment distance for a particular object in scene a to its supposed  
 1856 analogue in scene b. As a simple way to combine the edit costs across dimensions we first  
 1857 Z-score each dimension, such that the average distance between any two values across all  
 1858 objects and all scenes and dimensions is 1. We then take the L1-norm (or city block  
 1859 distance) as the cost for converting an object in scene (a) to an object in scene (b), or visa  
 1860 versa. Note this is sensitive the size of the adjustment, penalizing larger changes in  
 1861 position, orientation or size more severely than smaller changes, while changes in color are  
 1862 all considered equally large since color is taken as categorical. Note also that for  
 1863 orientation differences we also always assume the shortest distance around the circle.

1864 If scene (a) has an object that does not exist in scene (b) we assume a default  
 1865 adjustment penalty equal to the average divergence between two objects across all  
 1866 comparisons (3.57 in the current dataset). We do the same for any object that exists in (a)  
 1867 but not (b).

1868 Calculating the overall similarity between two scenes involves solving a mapping  
 1869 problem of identifying which objects in scene (a) are “the same” as those in scene (b). We  
 1870 resolve this “charitably”, by searching exhaustively for the mapping of objects in scene (a)  
 1871 to scene (b) that minimizes the total edit distance. Having selected this mapping, and  
 1872 computed the final edit distance including any costs for additional or removed objects, we  
 1873 divide by the number shared cones, so as to avoid the dissimilarities increasing with the  
 1874 number of objects involved.

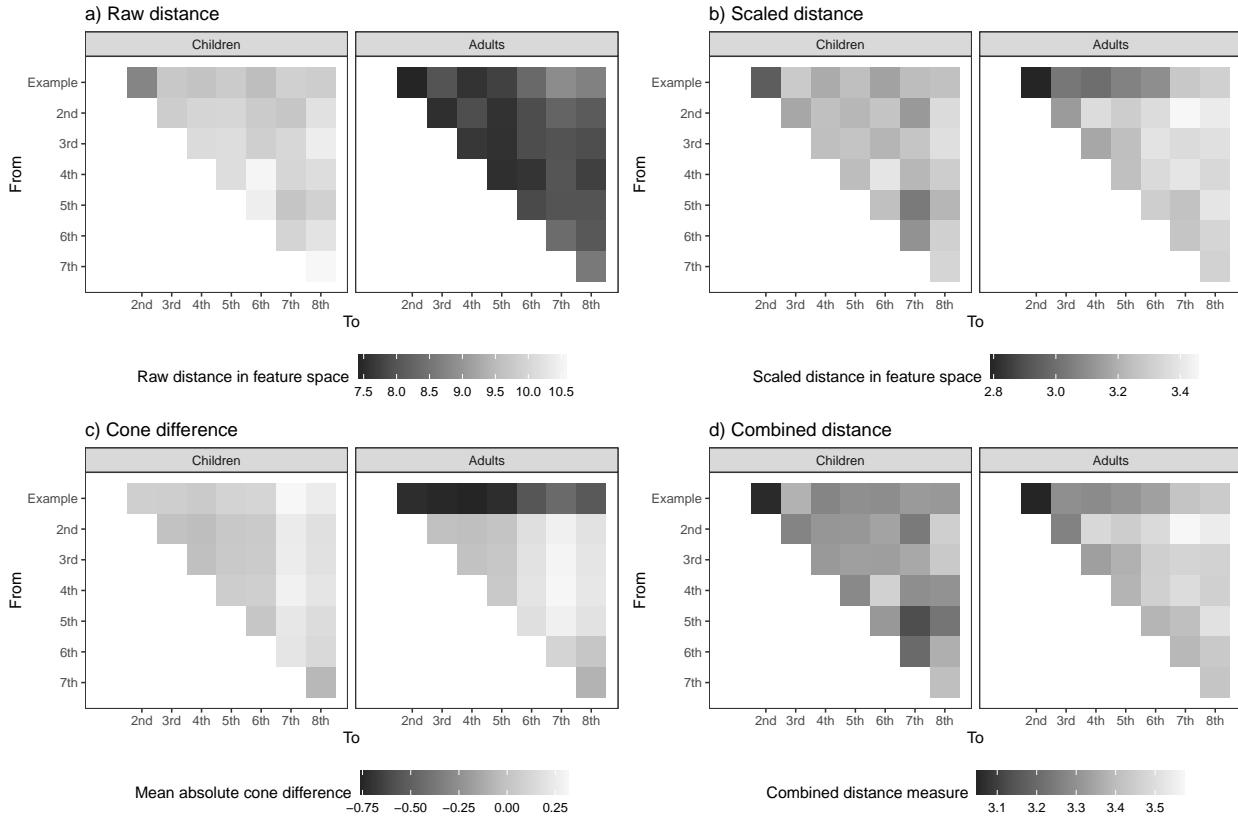
1875 Figure A-2 computes the inter-scene similarity components that go into Figure 6c in  
 1876 the main text. Summing up the edit distances across all objects, children’s scenes seem  
 1877 much more diverse than adults (Figure A-2a). However this is primarily due to their  
 1878 containing a greater average number of objects. Scaling the edit distance by the number of  
 1879 objects in the target scene gives a more balanced perspective (Figure A-2b) but does not  
 1880 account for the fact that the compared scene may contain more or fewer objects in total.  
 1881 Figure A-2c visualizes just the object difference showing that children’s scenes contain  
 1882 roughly as many objects on average as the initial example while adults’ scenes contain  
 1883 around 0.75 fewer objects than are present in the initial example (dark shading in top row).  
 1884 Thus, we opted to combine b and c by weighting the unsigned cone difference by the mean

**Figure A-1**

*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).*

<sup>1885</sup> inter-object distance across all comparisons to give our combined distance measure

<sup>1886</sup> (Figure A-2d and Figure 6c in the main text).

**Figure A-2**

*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.*

1887

#### Appendix D: Comparison with Bramley et al (2018)

1888

Finally, for interest and to demonstrate replication of our core results. We provide a

1889

direct comparison between the generalization accuracies in the current sample of children

1890

and adults and those in the sample of 30 adults modeled in (Bramley et al., 2018).

1891

Bramley et al (2018) included 10 ground truth concepts, and the current paper uses just

1892

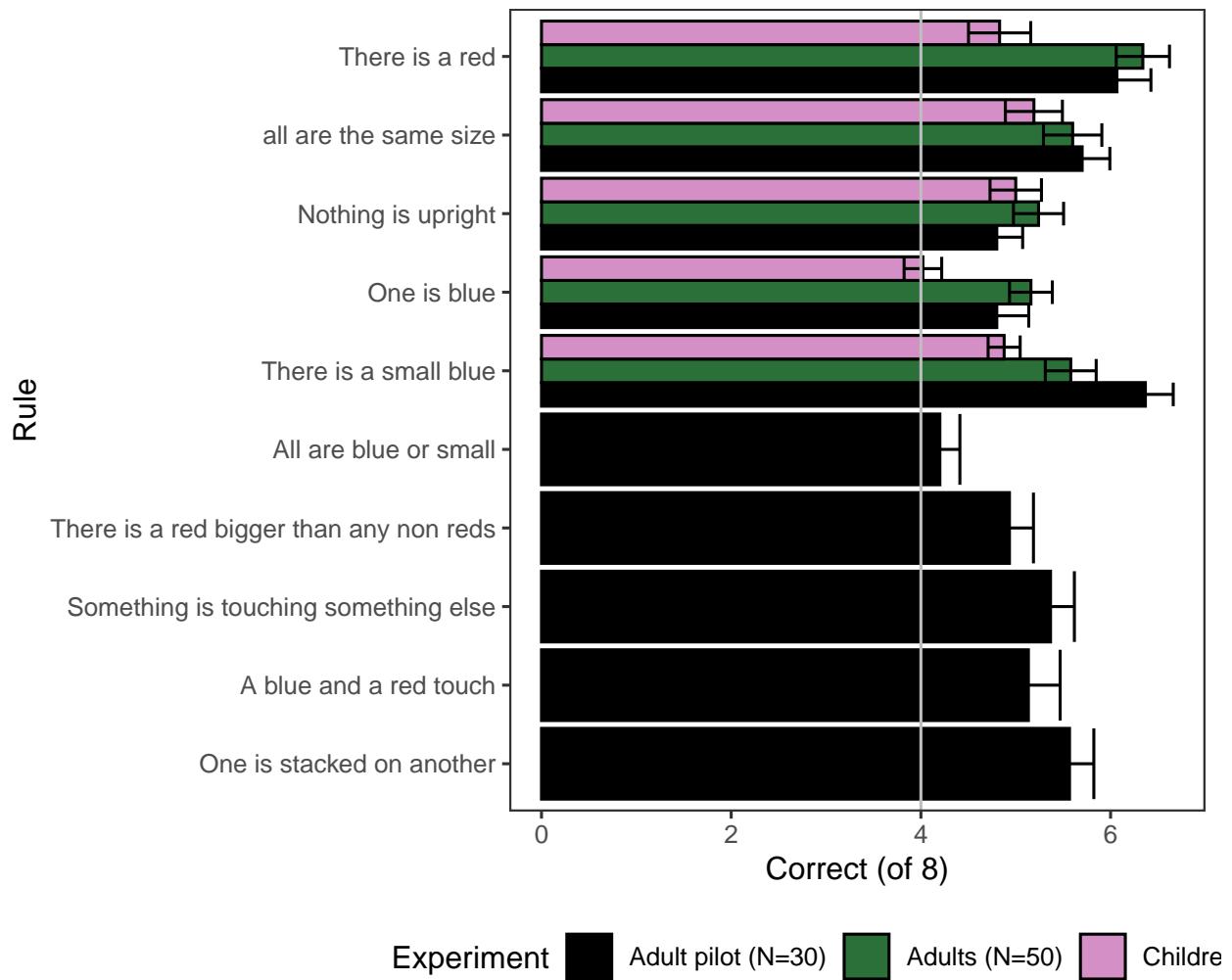
the first five of these. Figure A-3 shows these accuracy patterns side by side, revealing the

1893

adults in the current experiment performed approximately as well as those in the original

1894

conference paper.

**Figure A-3**

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