

Hypothesis generation through active inductive inference in children and adults

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. 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 ± 1.11) and 50 adults performed inductive inferences about a series of causal rules through active testing. Children generated substantially more complex guesses about the hidden rules and were more elaborate in their testing behavior. We take a ‘computational constructivism’ perspective to explaining these patterns, arguing that these inferences are driven by a combination of thinking (recombining and modifying existing symbolic concepts) and exploring (actively investigating and discovering patterns in the physical world). We show how our approach and rich new dataset speak to questions about developmental differences in hypothesis generation, active learning and inductive generalization.

Hypothesis generation through active inductive inference in children and adults

“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 new ideas and hypotheses to explain the world around them.
3 Children excel at forming new categories, concepts, and causal theories (Carey, 2009) and
4 by maturity, this coalesces into a capacity for intelligent thought characterized by its
5 domain generality and occasional moments of insight and innovation. Constructivism is an
6 influential perspective in developmental psychology (Carey, 2009; Piaget, 2013; Xu, 2019)
7 and philosophy of science (Fedyk & Xu, 2018; Quine, 1969) that posits learners actively
8 construct new ideas through a mixture of thinking—recombining and modifying ideas—and
9 play—exploring and discovering patterns in the world (Bruner, Jolly, & Sylva, 1976; Piaget
10 & Valsiner, 1930; Xu, 2019). In this way, constructivism views higher level cognition as
11 symbolic and at least somewhat language-like (Fodor, 1975) in its ability to make “infinite
12 use of finite means” (von Humboldt, 1863/1988). While the tenets and promise of
13 constructivist accounts are appealing, it has until recently lacked a formalization capable of
14 producing testable predictions in psychology experiments or detailed insights into human
15 cognition. In this paper we draw on recent methodological advances to formalize a
16 constructivist modeling framework and use it to analyze children and adults’ behavior in
17 an open-ended inductive learning task. Our account allows us to capture a wide range of
18 responses and behaviors and closely examine developmental differences in hypothesis
19 generation and active learning. To foreshadow, we argue that both children’s hypothesis
20 generation and active learning is based on a “flatter” generation process than adults,
21 producing greater diversity of ideas and behaviors but that is weaker in evaluative precision
22 and systematic coverage of the possibility space.

23 Concept learning

24 Classic work in experimental psychology suggests that some form of symbol
25 manipulation is required to synthesize the human capacity for reasoning and problem
26 solving (Bruner, Goodnow, & Austin, 1956; Johnson-Laird, 1983; Wason, 1968). However,
27 early symbolic accounts faced challenges in explaining how such representations could be
28 learned from sparse, noisy evidence, as well as how they could be effectively applied to
29 reasoning under uncertainty (Oaksford & Chater, 2007; Posner & Keele, 1968). Meanwhile,
30 various statistical accounts of concept learning have been developed that model concepts as
31 based in “family resemblance” within a feature space — for instance, centered around a
32 prototypical example or set of exemplars (Kruschke, 1992; Love, Medin, & Gureckis, 2004;

³³ Medin & Schaffer, 1978; Shepard & Chang, 1963). Such accounts have helped explain how
³⁴ people impute discrete categories from finite continuous evidence and how they then use
³⁵ these to generalize effectively to novel stimuli (Shepard, 1987). However, the core
³⁶ representations implied by this approach lack symbolic structure so are not able to explain
³⁷ the human capacity for generating conceptual novelty (Komatsu, 1992).

³⁸ More recently, Bayesian models have also played a major role in recent study of
³⁹ concept learning, providing a principled way of modeling probabilistic inference over both
⁴⁰ sub-symbolic and symbolic hypothesis spaces (Howson & Urbach, 2006). On the symbolic
⁴¹ side this includes inferences about particular causal structures (Bramley, Lagnado, &
⁴² Speekenbrink, 2015; Coenen, Rehder, & Gureckis, 2015; Gopnik et al., 2004; Steyvers,
⁴³ Tenenbaum, Wagenmakers, & Blum, 2003) as well as more general causal theories
⁴⁴ (Goodman, Ullman, & Tenenbaum, 2011; Griffiths & Tenenbaum, 2009; Kemp &
⁴⁵ Tenenbaum, 2009; Lucas & Griffiths, 2010). A key strength of the Bayesian framework is
⁴⁶ its account of how hierarchically structured representations support “learning to learn”—
⁴⁷ capturing the human capacity to make accurate few-shot inferences about familiar domains
⁴⁸ (Griffiths & Tenenbaum, 2009; Kemp, Goodman, & Tenenbaum, 2010). However, since
⁴⁹ Bayesian accounts, by definition, describe learning within a predefined hypothesis space,
⁵⁰ they are silent about how a learner might explore or generate possibilities within an infinite
⁵¹ latent space. For example, causal structure learning accounts using the Causal Bayesian
⁵² Network framework have typically assumed a fixed and finite—albeit sometimes very
⁵³ large—set of possible structural hypotheses (M. Jones & Love, 2011; Perfors, Tenenbaum,
⁵⁴ Griffiths, & Xu, 2011).¹ That is, hierarchical Bayesian accounts of induction are generally
⁵⁵ cast at Marr’s computational level (Marr, 1982). They show that people behave roughly *as if*
⁵⁶ they consider and average exhaustively over a weighted taxonomic hierarchy of
⁵⁷ possibilities and parameters, despite this being intractable in practice in nontrivial settings.

⁵⁸ Alongside Bayesian models, information theory has also featured frequently in
⁵⁹ cognitive science as a metric of idealized information acquisition (Gureckis & Markant,
⁶⁰ 2012), including choice of interventions and experiments that reveal causal structure
⁶¹ (Bramley, Dayan, Griffiths, & Lagnado, 2017; Bramley et al., 2015; Coenen et al., 2015;
⁶² Steyvers et al., 2003). However, information-theoretic analyses also presuppose the
⁶³ Bayesian notion that learners have the relevant possibilities in mind and act to
⁶⁴ discriminate between them, rather than to support the task of constructing or discovering
⁶⁵ better ones. Thus, while probabilities and information theory provide a jumping off point
⁶⁶ for top-down analysis of cognition, we should take their limitations seriously when seeking

¹ Although see (Buchanan, Tenenbaum, & Sobel, 2010) for an exception.

67 to reverse engineer humanlike inferential processing (Van Rooij, Blokpoel, Kwisthout, &
68 Wareham, 2019).

69 The central goal of this paper is to develop a computational model of constructivist
70 inference and use it as a lens to examine children’s and adults’ learning in a rich
71 open-ended task where the space of potential hypotheses and behaviors is effectively
72 unbounded. To achieve this, we focus on recent work in cognitive science that has
73 attempted to marry symbolic and statistical perspectives. This work characterizes
74 computational principles driving both human development and intelligence as resting on a
75 capacity to flexibly generate, adapt, combine and re-purpose symbolic representations when
76 learning and reasoning, but crucially doing so in ways that respect probabilistic principles
77 of inference under uncertainty (Bramley, Dayan, et al., 2017; Goodman, Tenenbaum,
78 Feldman, & Griffiths, 2008; Piantadosi, 2021; Piantadosi, Tenenbaum, & Goodman, 2016).

79 **Constructivism**

80 Fundamentally, we take the constructivist account to depart from
81 computational-level Bayesian accounts because it presumes representational *incompleteness*.
82 By this, we mean that the constructivist learner has not, and normally could not, consider
83 and weigh all the possibilities in play when learning. Instead, the constructivist is able to
84 *search* the hypothesis space and compare generated hypotheses to one another
85 (A. N. Sanborn & Chater, 2016; Stewart, Chater, & Brown, 2006). For example, this might
86 occur via processes that stochastically combine elements from an underlying concept
87 grammar to produce new variants or to adapt existing hypotheses in the light of new data
88 (Bonawitz, Denison, Gopnik, & Griffiths, 2014; Lewis, Perez, & Tenenbaum, 2014;
89 Nosofsky & Palmeri, 1998; Nosofsky, Palmeri, & McKinley, 1994). Outside of narrow
90 experimental settings, this modal incompleteness seems completely normal. A simple
91 illustration is the characteristic gap between generation of and evaluation of hypotheses or
92 examples (Gettys & Fisher, 1979). People have been shown to spontaneously generate only
93 a few explanations for a surprising event, such as a car engine’s failure, and fewer when put
94 under time pressure. However, they will endorse many more if a list of possibilities is
95 provided to them. Similarly, it takes a lot of time and effort to generate a list of category
96 members (such as all listing all fifty US States) while most Americans would find it
97 straightforward to identify the states among a list of state and non-states. Inference about
98 any area of active scientific inquiry, such as is reported in scientific journals like this one,
99 will typically involve an enormous latent space of potential explanatory theories only a
100 fraction of which have ever been articulated or tested. It is generally accepted that the
101 ground truth is unlikely to be among the set of theories already on the table (Box, 1976).

102 Fortunately, the idea that inductive inference involves symbolic search demystifies a
103 number of behavioral patterns that look like biases from the computational-level
104 perspective. These include order effects, anchoring, probability matching and confirmation
105 bias. Order effects are pervasive in human learning. If new hypotheses are arrived at
106 through a limited local search starting from a previous hypothesis then we should expect
107 patterns of dependency and autocorrelation between a single learner's hypotheses over time
108 (Bramley, Dayan, et al., 2017; Dasgupta, Schulz, & Gershman, 2016; Thaker, Tenenbaum,
109 & Gershman, 2017). Similarly, anchoring effects in estimation can be modeled as the result
110 of limited local adjustment from a salient starting point (Griffiths, Lieder, & Goodman,
111 2015; Lieder, Griffiths, Huys, & Goodman, 2017). Probability matching is another
112 phenomenon that is natural under a constructivist perspective. In experiments,
113 participants often choose options in proportion to their probability of being correct or
114 optimal rather than reliably selecting the best action, as we might expect if they had the
115 full posterior to hand (Shanks, Tunney, & McCarthy, 2002). However, probability
116 matching is actually often a *best case* scenario for a learner limited to using the the
117 endpoint of a local search as their guess (Bramley, Dayan, et al., 2017). Furthermore, it has
118 been shown that under a variety of plausible everyday utilities, a single-sample-based
119 decision can represent the best computation–accuracy tradeoff to a resource limited learner
120 (Vul, Goodman, Griffiths, & Tenenbaum, 2009). Confirmation bias is also pervasive in
121 human reasoning and active learning (Klayman & Ha, 1989) and hard to explain in purely
122 Bayesian terms. Wason (1960) famously asked participants to test and identify a hidden
123 rule and initially simply told them that the sequence 2–4–6 followed the rule. The intended
124 true rule was simply “ascending numbers” but participants frequently guessed more
125 complex rules such as “numbers increasing by two”. Analysis of participants’ tests revealed
126 that they frequently generated only tests expected to be rule-following under their
127 hypothesis (such as 6–8–12), so failing to adequately challenge and disconfirm their
128 hypothesis. On a constructivist perspective, learners can only base their exploration on
129 testing the hypotheses they have actually generated. To the extent that certain simple
130 hypotheses like “ascending numbers” were less likely to be generated on the basis of the
131 provided example (cf. Tenenbaum, 1999), it is not surprising that participants failed to
132 exclude these possibilities with their tests.

133 In the computational cognitive science literature, symbolic search ideas manifest
134 under the label of “learning as program induction”. Such models have begun to be applied
135 to synthesizing humanlike problem solving and planning and tool use (Allen, Smith, &
136 Tenenbaum, 2020; Ellis et al., 2020; Lai & Gershman, 2021; Lake, Ullman, Tenenbaum, &
137 Gershman, 2017; Ruis, Andreas, Baroni, Bouchacourt, & Lake, 2020; Rule, Schulz,

¹³⁸ Piantadosi, & Tenenbaum, 2018). We build on these in developing our account.

¹³⁹ **Accounts of Development**

¹⁴⁰ The “child as scientist” (Carey, 1985; Gopnik, 1996) — or more recently, “child as
¹⁴¹ hacker” (Rule, Tenenbaum, & Piantadosi, 2020) — perspective casts children’s cognition as
¹⁴² driven by broadly the same inferential processes as adults’ but at an earlier stage in a
¹⁴³ journey of construction and discovery. While children have been shown to be capable
¹⁴⁴ active learners (McCormack, Bramley, Frosch, Patrick, & Lagnado, 2016; Meng, Bramley,
¹⁴⁵ & Xu, 2018; Sobel & Kushnir, 2006), performing interventions or asking questions more
¹⁴⁶ informative or relevant to their learning than mere observations or random choices, there is
¹⁴⁷ also evidence that children’s ability to learn effectively from active learning data is more
¹⁴⁸ fragile than adults’. For example, Sobel and Kushnir (2006) found children were much less
¹⁴⁹ accurate at causal structure identification in “yoked” conditions — where they had to use
¹⁵⁰ evidence generated by someone else to learn — while adults are less effected, sometimes
¹⁵¹ able to learn about as well from others’ data as their own (Lagnado & Sloman, 2006). This
¹⁵² performance gap has been argued to stem from a mismatch between whatever idiosyncratic
¹⁵³ hypotheses are under consideration by the observer being different than those being tested
¹⁵⁴ by the active learner, making them less able to use the data to progress their theories
¹⁵⁵ (Markant & Gureckis, 2014). Relatedly, children’s active learning behavior is often more
¹⁵⁶ repetitive than adults’ (McCormack et al., 2016; Sim & Xu, 2017) and frequently more
¹⁵⁷ narrowly focused on testing a single hypothesis at a time (A. Jones, Bramley, Gureckis, &
¹⁵⁸ Ruggeri, in revision; Ruggeri & Lombrozo, 2014; Ruggeri, Lombrozo, Griffiths, & Xu,
¹⁵⁹ 2016). This might reflect a less developed working memory, restricting the number of
¹⁶⁰ hypotheses children can keep track of and compare to evidence. It has recently been shown
¹⁶¹ that children’s foraging behavior in large autocorrelated environments is driven more by
¹⁶² directed exploration and less by generalization than adults’ (Wu, Schulz, Speekenbrink,
¹⁶³ Nelson, & Meder, 2017). An early emphasis on exploration has been argued to be an
¹⁶⁴ effective solution to a lifelong explore–exploit tradeoff, since earlier discoveries can be
¹⁶⁵ exploited for longer (Gopnik, 2020). Constructivism also provides a potential explanation
¹⁶⁶ for transitions between developmental “stages”, characterized by occasional leaps forward
¹⁶⁷ in insight. For instance, Piantadosi, Tenenbaum, and Goodman (2012) demonstrate how a
¹⁶⁸ program induction account can reproduce a characteristic developmental phase transition
¹⁶⁹ in number understanding from grasping a few small digits to understanding the full real
¹⁷⁰ number line. Differences between childlike and adultlike behavior might also be captured
¹⁷¹ by parameterizable differences in search, potentially reflecting principles of stochastic
¹⁷² optimization (Lucas, Bridgers, Griffiths, & Gopnik, 2014). For instance, young children

have been found to be quick to make broad abductive generalizations from a small number of examples — e.g. readily imputing novel physical laws to explain surprising evidence (L. E. Schulz, Goodman, Tenenbaum, & Jenkins, 2008). Building on this finding, children's hypothesis generation and search has been framed as rationally "higher temperature" than adults' — producing more diversity of ideas at the cost of being noisier (Lucas et al., 2014). This is algorithmically sensible as optimization over high dimensional spaces is known to be more effective when proposals are initially large and decrease over time, as in *simulated annealing* (Van Laarhoven & Aarts, 1987). However, a high diversity of guesses might also simply reflect that children have a flatter latent prior than adults, with construction weights less tuned through experience, inherently producing more diverse hypotheses at the cost of producing high probability ones less frequently. A third possibility is that children's hypothesis generation might be driven more by *bottom-up* processing than adults'. With less established expectations to go on, children's hypotheses might more directly *describe* encountered patterns, while adults might rely more on their existing knowledge hierarchy to generate and constrain hypotheses in a *top-down* way (Clark, 2012). We will contrast children's and adults' hypothesis generation and active learning in a rich task setting that allows us to closely investigate and delineate between these ideas.

Task

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

continuous scene-specific $x \in (0,8)$, $y \in (0,6)$ positions and orientations $\in (0,2\pi)$. In addition to cones' individual properties, scenes also admit many relational properties arising from the relative features and arrangement of different cones. For instance, subsets of cones might share a feature value (i.e., be the same color, or have the same orientation) or be ordered on another (i.e., be larger than, or above) and pairs of cones might have relational properties like pointing at one another or touching. This results in an extremely rich implicit space of potential concepts.

We note that, by design, the dimensionality of this task makes it very difficult. As with Wason's 2-4-6 example, and genuine questions of scientific induction, the hard part of this task is not evaluating whether a candidate hypothesis can explain the data but rather generating the right hypothesis from among the infinitely many possibilities in the first place. As with the 2-4-6 task, there are always infinite data-consistent possibilities and many of these may be more salient than the ground truth. Without carefully chosen active testing, a learner will frequently fail to rule out simple possibilities that more parsimoniously capture the data than the ground truth, essentially being left with evidence that would not lead even an unbounded Bayesian agent to the correct answer.²

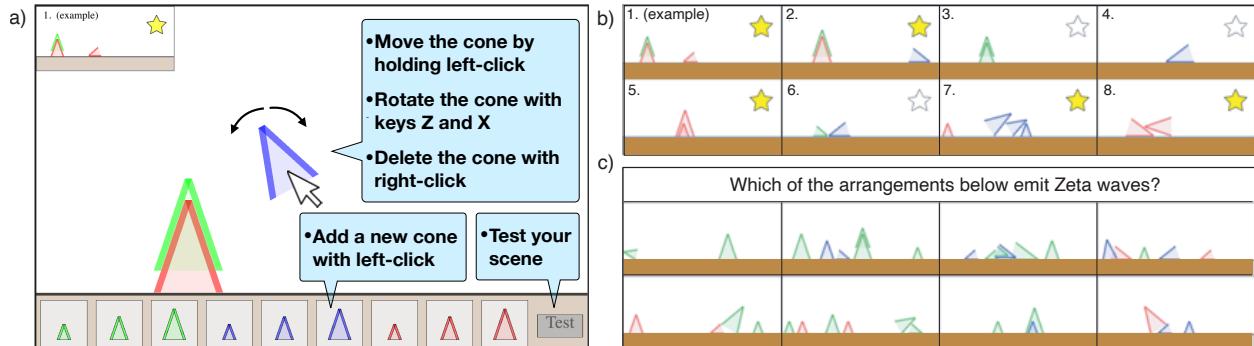


Figure 1

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.

We use mixed-methods (Johnson, Onwuegbuzie, & Turner, 2007), analyzing both qualitative data in the form of freely generated guesses about the symbolic rules and

² 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](#)).

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

234 Given the inherent difficulty of this type of task we expect absolute accuracy to be
 235 fairly low for both children and adults (and for our models). However, we expect that
 236 many participants will be able to make guesses that are consistent with most if not all of
 237 the evidence they have produced. Since we expect the evaluation of evidence–hypothesis
 238 consistency to be weaker in children, we expect adults’ guesses to be more strictly
 239 consistent with their evidence. While naively one might expect children to produce simpler
 240 guesses, our framework tends to predict the opposite. With a less trained construction
 241 probabilities, our models predict greater diversity and more elaborate structure in
 242 children’s guesses relative to adults’. While we will not attempt to model participants’
 243 scene generation process in detail, we expect similar principles to hold, with children’s
 244 scene generation being more baroque and elaborate than adults while simultaneously less
 245 closely tied to testing their hypotheses. Finally, there is the question of relative dominance
 246 of bottom-up and top-down processing in children’s and adults guesses. To explore this, we
 247 consider two models that vary in this dimension.

248 Context-free hypothesis generation

249 In accounting for children’s and adults’ inferences, we entertain two related
 250 constructivist algorithms. The first takes a fully “top down” approach to inference,
 251 utilizing a probabilistic context-free grammar (PCFG) to define a latent prior over concepts
 252 expressible in first order logic. A PCFG is a collection of “productions” that stochastically
 253 build expressions in an underlying grammar (Ginsburg, 1966). A PCFG can be used to
 254 generate a prior sample of hypotheses that can then be weighted by their likelihoods of
 255 producing observations—here, their ability to reproduce the labels for the scenes that the
 256 participant has tested. The model’s best guess about the hidden rule is then the *maximum*
 257 *a posteriori* hypothesis in the sample. The hypotheses make predictions about new scenes
 258 which can be weighted by their posterior probability and marginalized over to make
 259 generalizations. Because parts of this production process and underlying grammar involve
 260 branching—e.g., “and” and “or”—hypotheses can become arbitrarily long and complex,

261 involving multiple Boolean functions and complex relationships between an unlimited
 262 number of bound variables. In this way, an infinite latent space (in our case first order logic
 263 + lambda abstraction) is covered in the limit of infinite PCFG sampling (see Figure 2a).

264 The probabilities for each production in a PCFG can be fit to maximize
 265 correspondence with human judgments. Different PCFGs, containing different primitives
 266 and expansions, can be compared against human behavior. In this way, recent work has
 267 attempted to infer the “logical primitives of thought” (Goodman et al., 2008; Piantadosi et
 268 al., 2016). In the current work we consider a single expressive PCFG architecture but
 269 examine its behavior under uniform production weights but also with weights engineered to
 270 produce “childlike” and “adultlike” symbolic guesses. Crucially, under all three weighting
 271 schemes, our PCFG embodies the principle of parsimony: Simpler concepts—composed of
 272 fewer grammatical parts (Feldman, 2000)—have a higher probability of being produced and
 273 so are favored over more complex ones equally able to explain the data.

274 What PCFG approaches have in common is a generative mechanism for sampling
 275 from an infinite latent prior, here over possible logical concepts. However, sampled
 276 “guesses” must then be tested against data. Unfortunately, most samples are likely to be
 277 inconsistent with whatever data a learner has already encountered.³ For this reason, the
 278 procedure is inherently inefficient, and requires a very large numbers of samples in order to
 279 reliably identify non-trivial rules. Thus, we also consider an alternative that provide a
 280 more computationally plausible generation mechanism.

281 Context-based hypothesis generation

282 Instance Driven Generation (IDG) (Bramley et al., 2018) is a recent proposal
 283 related to the PCFG but with one key difference. Rather than generating hypotheses *a*
 284 *priori*, it generates ideas *inspired* by encountered examples (cf. Michalski, 1969), thus
 285 blending top-down generation with bottom-up reactivity to evidence. An IDG learner
 286 starts by observing the features of objects in a scene and uses these to back out a true
 287 logical statement about the scene in a stochastic but truth-preserving way. If the scene is
 288 rule following, this statement constitutes a positive hypothesis about the hidden rule.
 289 Otherwise, it constitutes a negative hypothesis, i.e. about what must *not* be present. Thus,
 290 IDG does not generate uniformly from all possible concepts, but directly from a restricted
 291 space consistent with a focal observation. Figure 2b illustrates this approach. While the

³ Many here are 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 distinct 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”).

292 PCFG starts at the outside and works inward, the IDG starts from the central content
293 (drawn from an observation) and works outward out to a quantified statement, ensuring at
294 each step that it is true of the scene. As with the PCFG, we consider a uniform variant as
295 well as variants that include productions reverse engineered to match the summary
296 statistics of guesses generated by children and by adults.

297 **Active learning**

298 Children have long been seen as primarily active learners, using “play” to explore
299 their environment and test their hypotheses (Bruner et al., 1976; Cook, Goodman, &
300 Schulz, 2011; Piaget & Valsiner, 1930). Information theory is used to benchmark active
301 learning (Nelson, 2005; Shannon, 1951) but assumes the learner knows the relevant
302 possibilities and acts to discriminate rather than to constructing or discover better ideas. It
303 has been argued that this may explain why behavioral alignment with information
304 maximization is mixed and task dependent (Coenen, Nelson, & Gureckis, 2018). In parallel
305 to rational accounts, the developmental literature has emphasized the utility of “control of
306 variables” heuristic (Chen & Klahr, 1999; Klahr, Fay, & Dunbar, 1993; Klahr, Zimmerman,
307 & Jirout, 2011) — this amounts to manipulating one design variable relative to a previous
308 test at a time to avoid confounded evidence, such that any changes in the outcome can be
309 unambiguously attributed to the change in the input. While not always the most efficient
310 strategy from a normative perspective (Coenen, Ruggeri, Bramley, & Gureckis, 2019;
311 A. Jones et al., in revision), a control of variables strategy intuitively minimizes the
312 cognitive costs of both testing hypotheses and of coming up with new scenarios (Gershman
313 & Niv, 2010). Past research has only focused on restricted settings with a few simple
314 variables and our task is much more complex. Nevertheless, in exploring the active learning
315 in our task, we will look for the empirical signature of control of variables in the form of
316 incremental and systematic testing patterns.

317 In sum, the core goal of this work is a close investigation of developmental
318 differences in active open-ended hypothesis generation examined through the lens of a
319 constructivist modeling approach that emphasizes the role of stochastic generation and
320 search. To foreshadow, we find children’s and adults’ guesses reflect a partially bottom-up
321 process of compositional concept formation. Children create more complex learning data
322 than adults but do so less systematically. They then go on to make more complex guesses
323 about the hidden rule that are only a marginally worse fit to evidence. Our constructivist
324 framework suggests this behavior is a natural result of “flat” idea and action generation
325 mechanisms. Crucially, our modelling then shows that both children’s and adults’ symbolic
326 guesses causally drive their generalizations, rather than these being driven by similarity

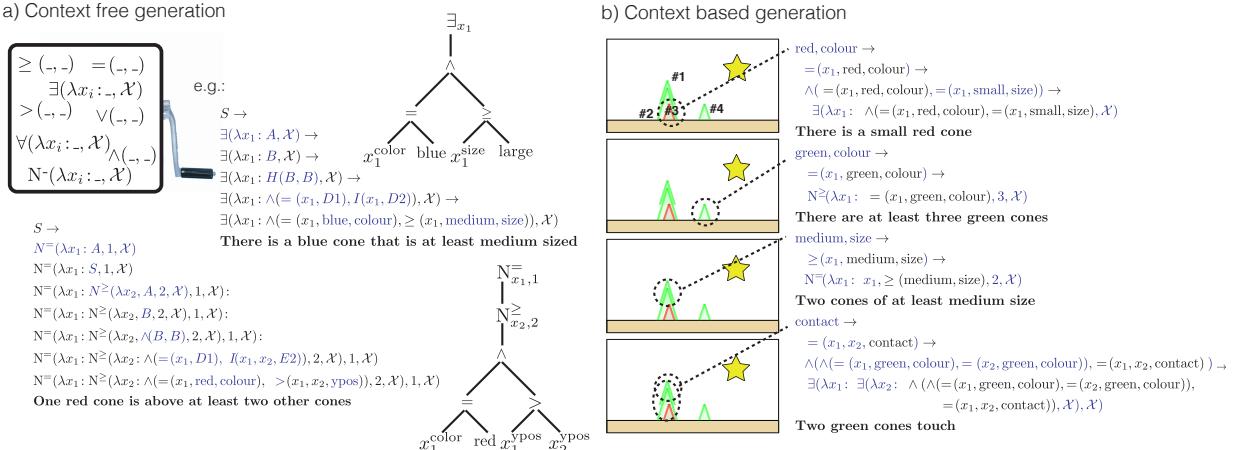


Figure 2

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.

³²⁷ and family resemblance.

Experiment

329 Methods

330 *Participants*

We recruited 54 children in the lab (23 female, aged 8.97 ± 1.11) and 50 adults online (22 female, aged 38.6 ± 10.2). Forty children completed all five trials and the remaining 14 completed 2.71 ± 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.⁴ 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 ± 0.75). Children's sessions lasted between 30 minutes and an hour. For adults, the task took 27.49 ± 12.09 minutes of which 9.8 ± 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.

⁴ 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.

Table 1
Rules Tested in Experiment

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

³⁴³ **Design**

³⁴⁴ All participants faced the same five learning problems in an independently
³⁴⁵ randomized order (see Table 1). For each learning problem participants were given an
³⁴⁶ initial positive example, as shown in the table, and then performed self tests of their own
³⁴⁷ before making generalizations and free guesses as to the hidden rule.

³⁴⁸ **Materials and Procedure**

³⁴⁹ **Child sample.**

³⁵⁰ **Instructions.** Participants sat in front of a laptop with a mouse attached, with
³⁵¹ the experimenter sitting next to them and interacted with the task through the browser.

³⁵² The experimenter read out the instructions for the participant. These explained
³⁵³ how the game worked and showed the participant five examples of possible rules the blocks
³⁵⁴ could have (relating to color, size, proximity, angle, or relation). The instructions also
³⁵⁵ included videos showing the participant how to manipulate the blocks using the mouse and
³⁵⁶ keyboard. After the instructions, the participant was given a comprehension check of five
³⁵⁷ true or false questions. If they did not get them all right on their first try, the experimenter
³⁵⁸ read through the instructions again and asked them again. All participants passed the
³⁵⁹ comprehension check the second time.

³⁶⁰ **Learning Phase.** The participant was then introduced to an initial example of a
³⁶¹ block type (“Here are some blocks called [name]s. We’re going to click test to see if stars

362 will come out of the [name]s.”). The initial example of each block type (i.e., each rule) was
363 constant across participants. Since every initial example of a block type was a positive
364 example, a star animation played when the “Test” button was clicked. The participant was
365 encouraged to use either the trackpad or the mouse to click the “Test” button, whichever
366 was comfortable for them.

367 After the initial positive example, the participant was shown a blank scene with
368 blocks available to add to it, and was asked to test the blocks seven more times
369 (Figure 1a). The scene creation interface was subject to simulated gravity, meaning there
370 were physical constraints on how the objects can be arranged. The experimenter told them
371 they could now play with the blocks like they saw in the instructional video. The
372 experimenter also reminded the participant of how to add, remove, move, and rotate blocks
373 on the screen using the mouse and keyboard. Participants were encouraged to ask for help
374 with moving the blocks if needed. If they seemed to be having trouble, the experimenter
375 would ask if they needed help with setting up the blocks. The participants were told that
376 when they had finished moving the blocks around, they should press the “Test” button to
377 see if stars came out of them. For positive tests, the experimenter would neutrally say:
378 “Stars did come out of the [name]s that time” and for negative tests: “Stars did not come
379 out of the [name]s that time.”

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

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

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

394 **Adult sample.** We recruited our adult sample from Amazon Mechanical Turk
395 and adults completed the task on their own computers. They completed the same
396 instructions as the children with an additional section about bonuses and had to
397 successfully answer comprehension questions, including an additional two about the

398 bonuses, before starting the main task. Specifically, adults were bonused 5 cents for each
 399 correct generalization (up to a possible 40 cents for each of the five trials) and an
 400 additional 40 cents for a correct guess as to the hidden rule, again for each of the five trials.
 401 Aside from having no experimenter in the room, and filling out the text fields themselves,
 402 the procedure was identical to the children’s task. Full materials including experiment
 403 demos, data and code are available at the [Online Repository](#).

404 Results

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

413 Rule complexity and constituents

414 We had human coders translate participants’ free text guesses about the hidden rule
 415 wherever possible into an equivalent logical lambda expression using the grammatical
 416 elements available to our learning models. We were able to do this for 86% (n=205) of
 417 children’s trials and 88% (n=219) of adults’ trials. For example, if the participant wrote
 418 “*There must be one big red block*” this was converted into
 419 $N = (\lambda x_1 : \wedge(=(x_1, \text{large}, \text{size}), =(x_1, \text{red}, \text{color})), 1, \mathcal{X})$. This logical version can be
 420 automatically evaluated on the scenes and can be read literally as asserting “*There exists
 exactly one x_1 in the set of objects \mathcal{X} such that x_1 has the size ‘large’ and the color ‘red’*”.
 421 We had a primary coder, blind to the experimental hypotheses code all responses, and a
 422 second coder blind spot check 15% of these (64). The two coders agreed in 95% of cases.
 423 We provide further details about the coding in Appendix B and full coding resources and
 424 full coding data in the [Online Repository](#).

426 To explore structural differences in children’s versus adults’ hypotheses, we first
 427 break down these encoded rule guesses into their logical parts. This reveals that children’s
 428 encoded rules were substantially *more complex* than those generated by adults and that
 429 both were more complex than the ground truth rules. Children’s and adults’ rules also
 430 differed in terms of the prevalence of particular elements and features (see Figure 3). As an
 431 example, one child’s rule for problem 1 was “*You must have two reds and one blue*” which

432 was translated to $N^=(\lambda x_1: N^=(\lambda x_2: (\wedge(= (x_1, \text{red}, \text{color}), = (x_2, \text{blue}, \text{color})), 1, \mathcal{X}), 2, \mathcal{X})$,
 433 requiring two quantifiers ($N^=$), one boolean (\wedge), 2 equalities ($=()$), and two references to
 434 the feature color. The typical child-generated-rule used 2.25 quantifiers (3b), 2.06 booleans
 435 (3c), 1.55 equalities and inequalities (3d), referred to 1.39 different primary features (color,
 436 size, orientation, x- or y-position, groundedness, 3e) and 0.37 relational features (contact,
 437 stackedness, pointing, or insideness, 3f). In contrast, the average adult generated rule
 438 required just 1.84 quantifiers, 1.20 booleans, 1.47 equalities and inequalities, and referred
 439 to 1.44 primary features but only 0.16 relational features. When children posited that an
 440 “at least”, “at most” or “exactly” a certain number of objects must have certain features,
 441 the number they chose were substantially higher than that for adults (2.36 compared to
 442 1.58). In terms of features, adults strongly tended to posit rules relating to color (58%
 443 compared to 39% of children’s rules), while children were more likely to refer to positional
 444 properties (26% compared to 18% of adults’ rules) and relations (31% compared to 14% of
 445 adults’ rules) between the objects.

446 Reverse engineering Childlike and Adultlike prior productions. Having
 447 encoded all the rule guesses from adults and children, we can work back from the
 448 distribution of rules to create a sets of productions that produce similar posterior samples.
 449 To do this, we work back from the observed counts for each rule element doing this
 450 separately for children and for adults. To roughly accommodate the fact that each guess is
 451 based on different learning data, we regularized these counts by including a prior
 452 pseudo-count of 5 on all productions. This value was not fit to the data, and simply serves
 453 to smooth the predictions a little. For example, children’s rules involved \exists 263 times, \forall 108
 454 times and N 297 times, so we assumed prior production weights of
 455 $\{263 + 5, 108 + 5, 297 + 5\}/(263 + 108 + 297 + 15) = \{.39, .17, .44\}$. A full set of fitted
 456 prior weights for both adults and children are visualized and detailed in Figure 3g&h.
 457 Strictly these are samples from a range of different participants’ posteriors $P(r|d_{p,t})$ not
 458 from their prior $P(r)$, since judgments were always conditional on some evidence. However,
 459 since evidence differs greatly across the rules we considered and scenes participants created,
 460 and since the structural elements of the grammar (Booleans, Quantifiers etc) are not
 461 tightly tied to scene-specifics, we feel this still provides an informative and useful
 462 elucidation of differences in a common set of productions that can produce children’s and
 463 adults’ hypotheses. This analysis illustrates that children’s construction process is “flatter”
 464 than adults’ under a constructive account, with a greater average entropy over the various
 465 production steps of this process 1.28 ± 0.50 bits compared to 1.03 ± 0.59 bits,
 466 $t(13) = 3.2, p = 0.007$. To avoid double counting the data in modeling subjects’ specific
 467 guesses, we also created separate agegroup-appropriate prior production weights for each

468 participant based on the guesses of the participants' from the same agegroup, holding out
 469 their own guesses, and used these to produce a unique held-out prior sample for each
 470 participant.

471 ***Accuracy***

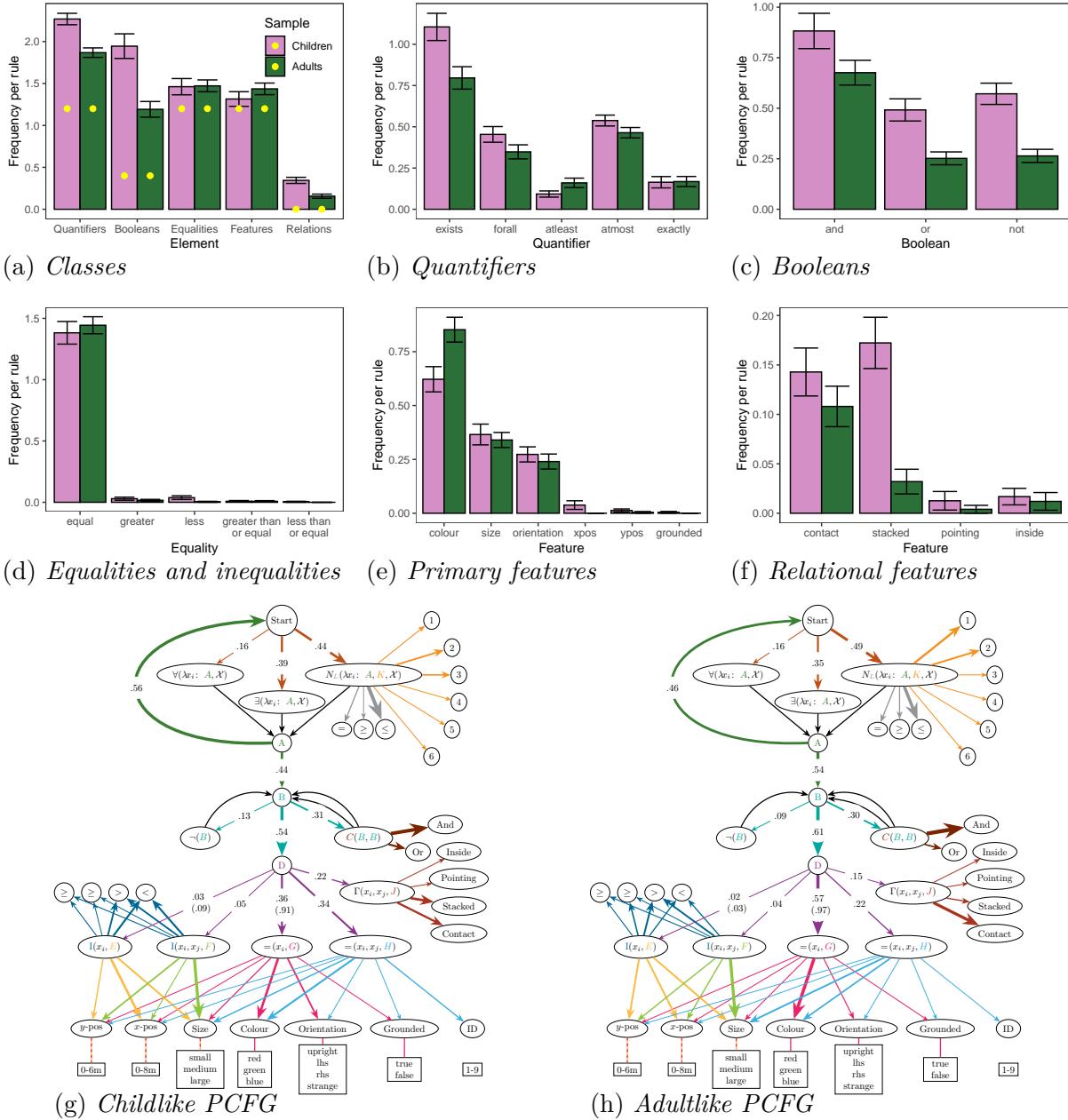
472 Having observed systematic differences in the content of children's and adults'
 473 hypotheses, we now ask if these manifest in children's and adults' inferential success; their
 474 ability to identify the ground truth and make accurate generalizations.

475 **Rule guesses.** Both children and adults were occasionally able to guess exactly
 476 the correct rules, doing so a respective 11% and 28% of trials. Adults produced the correct
 477 rule more frequently than children $t(102) = 4.0, p < .001$ and were more likely than children
 478 to guess correctly (at a corrected significance level of 0.01) for the "All are the same size",
 479 "One is blue" and "There is a small blue" rules (see Figure 4a). The plot reveals that no
 480 child identified rule 4 exactly "One is blue" and only one identified rule 5 "There is a small
 481 blue", while a slightly greater proportion of children than adults identified the positional
 482 "Nothing is upright" rule. Note that chance level baseline for these free guesses is
 483 essentially 0%. There are an unlimited number of wrong guesses and a small set of
 484 semantically correct guesses. It is also the nature of this inductive problem that there are
 485 an infinite number of wrong yet perfectly evidence-consistent rules for any evidence and
 486 often there is a simpler evidence-consistent rule available than the ground truth.⁵ Thus, it
 487 is instructive to ask whether participants' rules, where not exactly correct, are nevertheless
 488 consistent with the evidence they have generated.

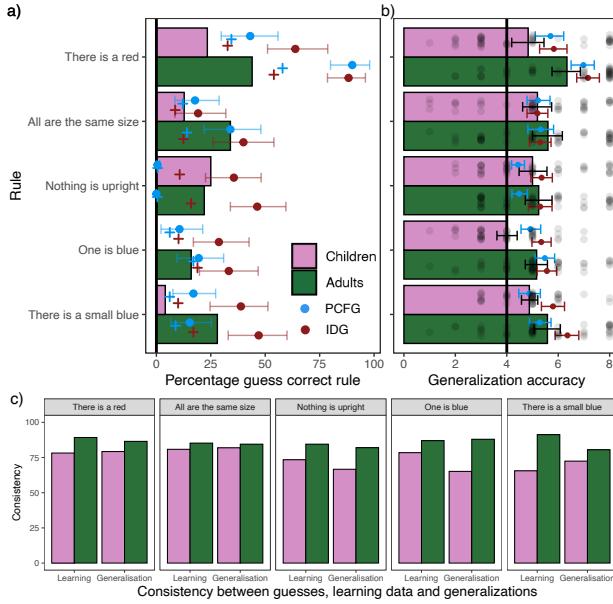
489 While, a completely random rule would only be consistent with all 8 scenes around
 490 $0.5^8 \times 100 = 0.4\%$ of the time, children's explicit rule guesses were perfectly consistent with
 491 the labels of the 8 training scenes 30% of the time and Adult's guesses were fully consistent
 492 54% of the time. There was a moderate difference in average proportion of the learning
 493 data explained by children's compared to adults' rules $71\% \pm 27\%$ vs $87\% \pm 17\%$
 494 $t(98) = 5.6, p < .001$. Similarly there was a difference the proportion of the participants'
 495 generalizations that were consistent with their rule guess $72\% \pm 21\%$ vs $84\% \pm 16\%$,
 496 $t(98) = 4.1, p < .001$ (see Figure 4c for a by-rule breakdown).

We now compare this to simulated *context free* (PCFG) and *context based* (IDG) hypothesis generation algorithms provided with the active learning data generated by the human participants. For each participant p — and separately for each learning task t in the case of the IDG — we generated 10,000 possible rules using *held-out*, age-consistent prior

⁵ Although as more evidence arrives the ground truth is increasingly likely to be among "simplest" rules in a posterior sample.

**Figure 3**

(a-f) Relative frequency of rule elements in Children's and Adults' rule guesses. Yellow points in a show ground truth frequency. (g&h) Visualization of the childlike and adultlike PCFGs (top down rule generation process), 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.

**Figure 4**

a) Percentage of participants guessing exactly the correct rule. Blue and red circles show performance of PCFG and IDG simulations guessing the maximum a posteriori rule, using age-group consistent prior production weights \pm bootstrapped 95% confidence intervals. “+” symbols shows the mean accuracy of a single posterior sample. b) Generalization performance. Bars show mean \pm bootstrapped 95% confidence intervals for children (pink) and adults (green). Model simulations select the maximum a posteriori label after marginalizing over the posterior. Black vertical lines denote chance performance. c) Consistency between subjects’ rule guess and their (self-generated) learning data, and generalization judgments.

productions derived above $\hat{R}_{\text{PCFG}_h}^p$ and $\hat{R}_{\text{IDG}_h}^{p,t}$ that have statistics matched to those in Figure 3a–f. We also generated an additional sample of 10,000 rules based on uniform production weights \hat{R}_{PCFG_u} , and similarly for the IDG generated a sample based on uniform productions for each task $\hat{R}_{\text{IDG}_u}^{p,t}$. The PCFG samples act as an approximation to the infinite latent prior over rules $P(r)$ before seeing any data, while the IDG samples are always generated conditional on one of the learning scenes. We split the IDG sample evenly across these 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 **d**, we then weight these samples by their likelihood of producing all eight scene

labels observed during the learning phase

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

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

and combine 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 500 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

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

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 1 for rules that predict all 8 scenes' labels to .13 for a single misprediction, and .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”.

Model rule guess accuracies are detailed in Table 2. Taking the *maximum a posteriori* (MAP) estimate (guessing in the event of ties) under all models leads to guessing the correct hypothesis with slightly higher frequency than participants (15–37% based on children’s active learning and 20–51% based on adults’). For instance, under a uniform-weighted prior sample, the PCFG MAP is correct on $15\% \pm 35\%$ of children’s trials and $20\% \pm 40\%$ of adults’ trials. Note that since these simulations use the same prior sample, the small differences we see are due to the different learning data generated by children and adults. However, accuracy improves substantially and better reproduces the empirical child–adult accuracy difference if we use samples based on reverse engineered weights that reproduce the qualitative properties of children’s and adults’ rules (see Appendix A and Figures 3g&h). For age-appropriate prior samples, the PCFG guesses

525 correctly on $18\% \pm 38\%$ of children’s trials and $32\% \pm 46\%$ of adults’ trials. Using an
 526 age-inappropriate “flipped” prior sample (i.e. child-like samples for adults and adult-like
 527 samples for children’s tasks) obliterates this difference $23\% \pm 42\%$ for children’s trials and
 528 $22\% \pm 41\%$ for adults’ trials. We see a similar pattern for the IDG algorithm, but higher
 529 accuracy across the board. The IDG has the better accuracy on both children’s and adults’
 530 trials, guessing over half of the hidden rules correctly ($51\% \pm 50\%$) in the case of adults’
 531 trials. However, achieving this accuracy requires maximizing over the full sample, while we
 532 have argued that process level accounts are more likely to yield behavior closer to posterior
 533 sampling (Table 2, right hand columns). Indeed posterior samples provide a visually closer
 534 fit to the by-rule guess rates (Figure 4a). To check what provides the best account of
 535 participants accuracy patterns we fit logistic mixed-effect regression model using each
 536 algorithm \times prior combination to predict participants’ by-task probability of guessing
 537 correctly, including random effects for both rule type and participant. The “Fit” columns
 538 of Table 2 shows the log likelihood for each of these models, revealing that participants’
 539 correct judgments were best predicted by posterior sampling under Instance Driven
 540 Generation with an age-appropriate prior (log likelihood = 211.5,
 541 $\beta = 5.44 \pm 1.74$, $Z = 5.99$, $p < .001$) improving over a baseline of -234.3 for a model with
 542 only intercept and random effects.

Table 2
Accuracy of Rule Guesses by Simulation Models

Algorithm	Prior	Acc Max Posterior (%)			Acc Sample Posterior (%)		
		Children	Adults	Fit	Children	Adults	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
IDG	Agegroup	37 ± 48	51 ± 50	-229	14 ± 16	24 ± 22	-216
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 participants’ correct responses.

543 **Generalizations.** We now use the models to account for the quantitative response
 544 data constituted by forced choice generalizations about which of 8 new scenes will produce
 545 stars (i.e. follow the hidden rule). Across the five tasks, both children and adults guessed
 546 more accurately than chance (50%): *children* $mean \pm SD$ $59\% \pm 11\%$, $t(53) = 5.9$, $p < .001$;
 547 *adults* $70\% \pm 14\%$, $t(49) = 10.3$, $p < .001$. Adults’ generalizations were significantly more

548 accurate than children's $t(102) = 4.6, p < .001$ and children's accuracy improved
 549 significantly with age $F(1, 52) = 6.2, \eta^2 = .11, p = 0.015$. Indeed, adults' generalization
 550 accuracy was above a Bonferroni-corrected chance level of $p \leq 0.01$ for all five rules and
 551 children were similarly above chance except for rules 1. "There is a red"
 552 ($t(46) = 2.5, p = .015$) and 4. "One is blue" ($t(46) = .1, p = .915$; see Figure 4b).

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

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

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

553 Provided with the active learning data generated by the human participants, both
 554 performed in the human range at generalization. Using uniform production weights and
 555 taking the marginally most likely generalization labels over a posterior weighted sample of
 556 PCFG-generated rules based on the participants active learning data yielded accuracies of
 557 $61.4\% \pm 19.6\%$ for children's and $63.5\% \pm 20\%$ for adults' data. Using agegroup-appropriate
 558 priors, PCFG model performs a little better and reproduces the empirical difference
 559 between children's and adults' accuracy: $62.8\% \pm 19.8\%$ for children's and $68.8\% \pm 20.9\%$ for
 560 adults' PCFG weights. The uniform IDG, again, performed slightly better than the PCFG,
 561 generalizing at $65.2\% \pm 19.3\%$ from children's active learning data and $69.0\% \pm 21.0\%$ from
 562 adults'. Again using agegroup-appropriate prior productions, the IDG models' performance
 563 increased slightly for both adults' and children's data and better reproduced the difference
 564 between children's and adults' accuracy ($68.8\% \pm 20.1\%$ and $74.2\% \pm 21.7\%$).

565 The stronger generalizations of the IDG compared to the PCFG replicates the
 566 findings of Bramley et al. (2018) and extends this to children as well as adults. Intuitively,
 567 this is because the bottom-up mechanism ties the hypotheses generated to features of the
 568 learning cases, effectively narrowing in on plausible hypotheses more efficiently. More
 569 broadly, these simulation results underscore the inherent difficulty of this task in particular
 570 and open-ended inductive inference in general. The PCFG and IDG were not statistically
 571 better or worse than participants at any rule inference under after Bonferroni correction
 572 with the exception that the IDG outperformed children on rule 4 $t(96) = 4.7, p < .0001$.
 573 Thus strikingly, even in this "small world" with known and fully observed features, and
 574 even allowing simulations to sample and maximize over implausibly large numbers of
 575 hypotheses, we could not robustly outperform human adults in this task. Additionally, we
 576 saw that building in human inductive biases boosted performance and that adults' stronger

577 inductive biases go some way to explain differences in generalizations.

578 ***Interim discussion***

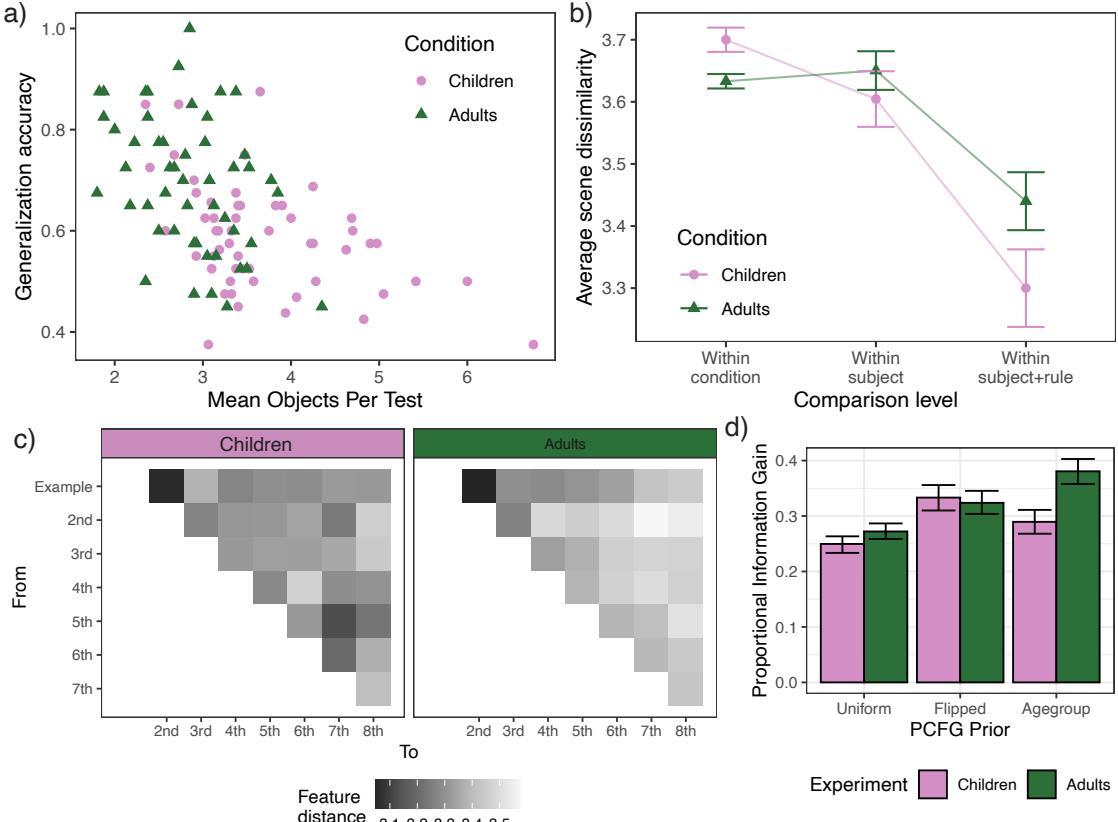
579 Children were only moderately less able to guess rules that fit the evidence than
 580 adults and there were only moderate differences in the compatibility between children's
 581 and adults' rules and their generalizations. However, children appeared to overfit the
 582 evidence more, essentially producing more complex and naïve characterizations of the
 583 rule-following scenes than did adults. This can be seen in the larger number of quantifiers
 584 and relations mentioned in children's rules than in adults', essentially referring to complex
 585 properties of the learning scenes that were irrelevant to their label.

586 A complicating factor is that children generated different learning data to adults.
 587 However, our PCFG and IDG simulations suggest exposure to different data cannot explain
 588 most of the accuracy differences between children and adults. Using identical production
 589 weights and the scenes generated by adults and children led to only small differences in
 590 accuracy for the PCFG and moderate for the IDG, while using a “flatter” set of productions
 591 fit to match childlike rules, and a more “peaked” set fit to adults' rules, better reproduces
 592 the accuracy differences. We take this to suggest hypothesis construction differences drive
 593 a large portion of the differences in children's and adult's inductive inferences.

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

597 ***Search behavior***

598 As well as generating more complex rules, children also tended to create more
 599 complex scenes than adults during the learning phase. The average child-generated scene
 600 contained 3.7 ± 0.88 objects (about the same as were present on average in the initial
 601 examples) compared to 2.8 ± 0.57 objects for adults ($t(102) = 5.8, p < .001$). The
 602 complexity of test scenes was inversely related to performance overall
 603 ($F(1, 102) = 39.0, \beta = -0.08, \eta^2 = .28, p < .001$) and also within both the child sample
 604 ($F(1, 52) = , \beta = -0.056, \eta^2 = .20, p < .001$) and adult sample
 605 ($F(1, 49) = 9.1, \beta = -0.096, \eta^2 = .16, p < .001$) taken individually (see Figure 5a). Within
 606 the child sample, age was inversely associated with scene complexity with an average of 0.35
 607 fewer objects per scene for each additional year of age $F(1, 52) = 12.6, \eta^2 = .19, p < .001$.
 608 Aside from this difference, we can also assess whether children's or adults' scenes bear the
 609 hallmarks of a local “search” across possible scene dimensions.

**Figure 5**

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

610 **Scene sequences and similarity.** While we do not yet have a model of scene
 611 creation process, we hypothesized that *control of variables* (Kuhn & Brannock, 1977) is a
 612 reasonable marker of systematic active learning. In the current setting, this manifests as a
 613 tendency to generate new evidence that comes close to recreating a previous scene (i.e.
 614 whose labels is already known) but making some change to it. This allows a learner to
 615 isolate boundary conditions for membership of the rule-following or non-rule-following
 616 scenes, and so potentially fine tune a focal hypothesis, or choose among a small set of
 617 similar alternatives (Bramley, Dayan, et al., 2017). Additionally, we speculated that reuse
 618 in general is likely to reduce cognitive load (Gershman & Niv, 2010).

619 If this is the case, we should expect the scenes generated by participants to be more
 620 similar to the initial example than to a random scene or a scene drawn from a different
 621 learning problem. To explore this, we constructed a distance metric that we used to

measure the featural 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 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 ± 0.14 compared to 3.63 ± 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 ± 0.459 , adults: 3.44 ± 0.33 , $t(102) = 1.77, p = 0.08$, Figure 5b&c).

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

Children’s and adults’ scene generation patterns manifest in differences in the quality of the total evidence generated according to an information gain analysis. For example, using the unweighted PCFG sample, prior entropy is 7.74 bits and children’s evidence produces an information gain (reduction in uncertainty) of 1.93 ± 0.45 bits while adults’ data average an information gain of 2.11 ± 0.38 bits $t(102) = 2.12, p = 0.035$ (see Figure 5d). Relative to the fitted PCFG priors, the difference in information gains is rather larger, with children’s scenes leading to information gain at 2.28 ± 0.66 bits (prior entropy 7.87 ± 0.05), and adults’ at 2.96 ± 0.64 (prior entropy 7.77 ± 0.04) $t(102) = 5.3, p < .0001$. Under the flipped priors — that is, taking the adultlike PCFG prior for children and childlike PCFG prior for adults — children’s tests look more informative than under their

658 own prior, generating 2.58 ± 0.68 bits, and adults' tests slightly less informative than under
659 their own prior 2.55 ± 0.57 bits, eliminating the statistical difference $t(102) = 0.24, p = 0.81$.
660 On the face of it, this is evidence against the idea that children's more elaborate hypothesis
661 generation and concomitantly flatter construction weights are driving them rationally
662 toward more elaborate testing choices. However, we see this information-theoretic analyses
663 as limited in what reveals. This is because is predicated on an implausibly complete
664 representation of uncertainty that we approximated by using a large sample of prior
665 hypotheses, while we might expect constructivist search behavior to be driven by more
666 focal testing of a smaller number of possibilities. Nevertheless, we present these information
667 scores as norms for completeness and comparison with other active learning tasks.

668 Experiment Discussion

669 Our analysis through the constructivist lens suggests we cannot attribute the
670 differences in children and adults' hypothesis generation to their differences in active
671 learning, nor can we attribute their differences in active learning behavior directly to
672 differences in hypothesis generation. That is, assuming the same hypothesis generation
673 process for children and adults' data does not reproduce the differences in rule guesses and
674 accuracy. At the same time, assuming children and adults' scene creation is driven by
675 distinguishing among *different* sets of hypotheses, sampled from a childlike or adultlike
676 latent prior, does not explain the developmental differences in the complexity and
677 systematicity of the scene creation. Rather, these data support the idea that
678 developmental shift in hypothesis generation and active search are manifestations of a
679 gradual sharpening of the constructivist generative mechanisms that produce relevant ideas
680 but also novelty in cognition.

681 We finally turn to a model-based analysis of the free responses and generalizations.
682 To foreshadow, we find both children's and adults' guesses are better accounted for by our
683 partially bottom-up IDG account of hypothesis generation than by our fully top-down
684 PCFG norm. We then find that both children's and adults' generalizations cannot be
685 explained by a non-symbolic "family resemblance" model but are well predicted by their
686 individual symbolic rule guesses.

687 Model comparison

688 Guesses

689 To evaluate our constructivist PCFG and IDG models' ability to explain
690 participants' free response guesses, we estimated probability of each approach generating

691 exactly the participant’s encoded guess based on their active learning data.

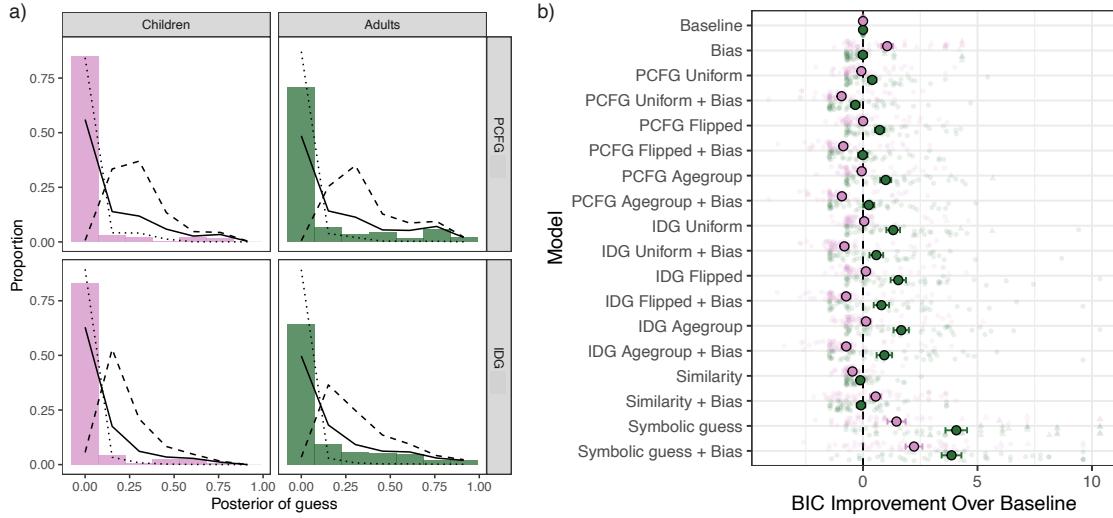
692 By definition, all 87% of participants’ rules that we were able to encode in our
 693 concept grammar have nonzero support under a PCFG prior, and due to the stochasticity
 694 we assumed in our likelihood function, all possibilities should also nonzero have posterior
 695 probability. However, in practice it is impossible to cover an infinite space of discrete
 696 possibilities with a finite set of samples, meaning there are a substantial number of cases in
 697 which we did not generate the participants’ guess. The proportion of Children’s and
 698 Adults’ rules that were generated at least once in 10,000 samples was 45% for children’s
 699 and 50% for adults’ rules from the uniform prior and 42% for children’s and 53% for
 700 adults’ guesses from the age-appropriate prior. The IDG samples included children’s rules
 701 49% of the time and adults’ 62% of the time using uniform weights and using fitted weights
 702 this increased to 69% of children’s and 76% of adults’ rules.

Table 3*Models Chance of Producing Participants’ 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
IDG	Agegroup	4.5 ± 7.1	15	14.1 ± 13.6	22

703 Table 3 details model fits to participants’ guesses. This suggests the IDG is the
 704 stronger hypothesis generation candidate, assigning higher probabilities on average to the
 705 rules that participants guessed. As expected, the variants of the PCFG and IDG with
 706 fitted production weights are better aligned with participants’ guesses than variants with
 707 uniform (or mismatched) weights. However, all models produced adults’ guesses with a
 708 much higher probability than children’s guesses.

709 Figure 6a additionally visualizes participants guesses in terms of their posterior
 710 probability under PCFG and IDG sampling and compares this to what we would expect if
 711 guesses are samples from the posterior (black line), the result of finding the maximum a
 712 posteriori guess of the 10,000 considered hypotheses (dashed line) or else are simply
 713 samples from the prior (dotted line). This visualization shows that, under all the models
 714 we consider, adults’ guesses are distributionally more consistent with posterior sampling
 715 than posterior maximization, while children’s appear somewhere between prior and
 716 posterior sampling.

**Figure 6**

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.

717 Generalizations

718 A standard benchmark for models of concept learning is a fit with participants'
 719 generalizations to new exemplars. Thus, we complete our analyses by comparing a range of
 720 models' ability to account for participant's specific generalizations. The set of models we
 721 consider allows us to test our core claims that children's and adults' induced representations
 722 are symbolic and compositional, as opposed to statistical and similarity-driven.

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

728 The models we fit were:

729 **1. Baseline.** Simply assigns a likelihood of .5 to each generalization $\in \{\text{rule}$
 730 $\text{following}, \text{not rule following}\}$ for each of the 8 generalization probes for each of the 5
 731 learning trials.

732 **2. Bias.** Acts a stronger baseline by allowing participants to have an overall bias

toward or against selecting generalization scenes as rule following. For this model, $b = 1$ if >50% of generalizations predict the scene is rule following and 0 otherwise. The model is fit using a mixture parameter λ to mix this modal prediction with the baseline prediction of .5 $P(\text{choice}) = \lambda b + (1 - \lambda).5$.

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

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

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

15-16. Similarity {No Bias, Bias}. Inspired by Tversky’s statistical and similarity based *contrast model of categorization* (cf., Tversky, 1977), we used the inter-scene similarity between each generalization scene and each training scene to compute the relative average similarity of each generalization case to the rule-following vs. the not rule-following training scenes. Similarities were computed 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 the inverse logit of these scores $r = \frac{\Delta\text{Similarity}}{\Delta\text{Similarity} + 1}$ 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 $r \in R : \{rule-following=1, not\ rule-following=0\}$ 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 PFG and IDG models to also fit generalization patterns even though they are blind to 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 6b.

Results

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 with a larger BIC improvement of 996 over Baseline, with a $\lambda = 0.08$ indicating less bias and temperature $\tau = 0.50$ indicating tighter alignment with the guessed-rule's predictions. Probing this bias, we see children undergeneralized substantially on average, selecting just $2.75 \pm 1.42/8$ scenes compared to adults' $3.42 \pm 1.03/8$ (unknown to the participants, there were always 4 rule following generalization scenes). Focusing on individual fits, the picture is mixed for children's generalizations, with 16/50 best fit by the *Bias* only model, followed

800 by 15 by the *Symbolic guess* model, 9 by the *Symbolic Guess + Bias* model and a further 7
 801 by the fully random *Baseline*. No other model best fit more than 2 children. For adults,
 802 32/52 were best fit by *Symbolic guess*, 6 by *Bias*, 4 by *Symbolic guess + Bias* and no other
 803 model best fit more than 2 participants. Overall, children’s generalizations are much harder
 804 to predict than adults’ with end-to-end constructivist accounts of their generalizations
 805 performing close to *Baseline*. This is partly to be expected since our child-like construction
 806 weights inherently produce a very diverse set of guesses and correspondingly diffuse set of
 807 generalization predictions. However, conditioning on Children’s symbolic guesses we were
 808 able to predict their generalizations far better than by *Similarity*, *Bias* or any other model
 809 we considered. Adults’ generalizations seem more straightforwardly driven by their
 810 symbolic guesses, with better individual fits on average using their guess directly without
 811 adjusting by any bias toward or against predicting scenes to be rule-following. This makes
 812 sense since, with a clear hypothesis in mind, there is little rationale to select more or fewer
 813 than than the generalization scenes consistent with that rule.

814 Thus, confirming our key constructivist hypothesis, there is a clear alignment
 815 between participant’s symbolic rule guesses and their generalizations. Thus, these results
 816 are consistent with our computational constructivism framework, with generalization
 817 behavior best accounted for as driven by the logical structure symbolic hypotheses rather
 818 than family resemblance. Our account explains why inference patterns in this setting are
 819 highly idiosyncratic to different participants yet explicable as the result of a process of
 820 compositional of construction and search. As with the free rule guesses, the IDG was
 821 better aligned with participants’ generalizations than the PCFG, particularly for adults,
 822 and particularly when using Agegroup-specific rather than Uniform or Flipped production
 823 weights. Thus, our model fitting also supports the idea that participants were inspired by
 824 patterns present in the learning data, such as the objects and relations in the initial
 825 positive example. However, this does not appear to be a developmental difference per se,
 826 with both children’s and adults’ judgments better accounted for by our IDG than our
 827 PCFG algorithm across all of our analyses.

828 Discussion

829 We explored children and adults’ active hypothesis generation and inductive
 830 inference in an interactive task where the space of possibilities and actions is
 831 compositional, open and practically unbounded. Our results are rich and nuanced but
 832 broadly we found that:

- 833 1. A constructivist framework could explain the diversity and distribution of people’s
 834 symbolic hypotheses in our inductive learning task, synthesizing both their sporadic

835 correct guesses but also their incorrect ideas and also offering an explanatory
836 framework for exploring differences between children’s and adults’ generation
837 mechanisms.

- 838 2. Children’s guesses and active testing reflect “flatter” generation mechanisms than
839 adults’, producing more diversity but less predictability.
- 840 3. The logical form of both children and adults’ symbolic guesses predicts their
841 generalizations to new scenes far better than feature similarity.
- 842 4. We found evidence of probability matching, with the frequency adults made
843 particular guesses more similar to posterior probability than the pattern expected
844 under noisy selection of the most probable hypothesis.
- 845 5. Both children and adults’ hypothesis generation was context sensitive, as shown by
846 better fit throughout our analyses for a partially bottom-up Instance Driven
847 generation account — sensitive to patterns in the learning data — over a fully
848 top-down Context Free (PCFG) generation account — aware only of the relevant
849 features and possible values.

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

853 **Limitations**

854 ***Experimental Control***

855 While this task and new dataset provide an exceptionally rich window on
856 developmental differences in inductive inference, some of what is gained in open-endedness
857 is lost in experimental control. There is residual ambiguity about the extent that
858 differences in active learning caused differences in hypothesis generation and visa versa.
859 Partialing this out would require experiments that fix the evidence and probe the
860 hypotheses generated, or that fix the hypotheses in play and probe what evidence is sought.
861 However, we have argued that such constrained tasks run the risk of short-circuiting
862 natural cognition. Learners may struggle to test hypotheses they did not conceive
863 themselves, and are known to struggle to use data they have not generated to evaluate their
864 hypotheses (Markant & Gureckis, 2014; Sobel & Kushnir, 2006). We take this to mean that
865 sole quantitative focus on studies that fix one or other aspect of the problem may provide a
866 misleading perspective on end-to-end active inference in the wild. We feel that our open
867 ended task provides a valuable complementary perspective to more constrained tasks.

868 ***Theoretical Expressivity***

869 There are many ways we could have set up the primitives, parameters and
870 productions of our PCFG and IDG. This makes for a dangerously expressive set of theories
871 of cognition. We do not claim to have explored this space exhaustively here but rather that
872 our modeling lends support to the idea that some symbolic and compositional process
873 drives children and adults' active inductive inferences about the world. That is, we can
874 explain the variability and productivity of human hypothesis belief formation in symbolic
875 terms. Identifying the computational primitives of thought may not be a realistic empirical
876 goal since a feature of constructivist accounts is their flexibility. Learners can grow their
877 concept grammar over time, caching new primitives that prove useful (Piantadosi, 2021).
878 Moreover, it is well known many different symbol systems can mimic one another (Turing,
879 1937), meaning that expressivity alone cannot distinguish between them. Since, we expect
880 different learners to take different paths in an inherently stochastic learning trajectory, this
881 limits universal claims about representational content.

882 ***Feature selection***

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

890 ***Sampling differences between children and adults***

891 One potential concern is that the complexity of children's guesses relative to adults
892 stems partly from their being collected verbally and in the presence of an experimenter
893 rather than typed during an online experiment. Speaking carries different cognitive
894 demands than typing and may lead to children simply responding in a more verbose way
895 than adults. While we cannot rule this out, we do not think this is a major concern.
896 Adults were well compensated for accuracy, meaning their motivation was primarily to be
897 correct rather than brief. The semantic content of both children's and adults rules were
898 extracted through our coding of them into lambda calculus meaning that surface
899 differences in concise expression can be separated from logical complexity. Furthermore
900 children's guesses were not the only thing that was more elaborate about their behavior.
901 They were also more elaborate in their active testing choices, producing more complex

902 scenes despite having to create these in the same manner as adults. Since the testing
903 interface was reset on each trial, this complexity took more effort, with children's scenes
904 requiring substantially more clicks and more time to produce than adults'.

905 ***Use of verbal protocols***

906 Another worry about our use of free responses is that they rely on a capacity for
907 precise linguistic expression not to mention the assumption that learners have insight into
908 the structure of their own concepts. It is known that children's vocabularies differ from
909 adults', raising the concern that some of our results reflect language use rather than the
910 concepts being articulated. While our artificial environment contains only simple objects
911 and basic features that are familiar to even young children, there is evidence that children's
912 speech does not distinguish as well among quantifier usage (e.g., all, each, every) until late
913 in childhood (Brooks & Braine, 1996; Inhelder & Piaget, 1958). Thus, it could be that
914 linguistic imprecision is behind some of the differences between children's and adults'
915 guesses. For instance, this seems like a potential explanation for the lack of any exactly
916 correct guesses from children about the quantifier-dependent rule 4 "exactly one is blue".
917 However, a closer look at responses reveals that only 11/47 children guessed a rule that
918 mentioned blue at all. Meanwhile 37/50 of adults' rules mentioned blue, but all but seven
919 of these were wrong about the particulars of the quantification. In many cases other
920 potential quantifications were not ruled out by adults' testing. For instance, several
921 subjects never tried adding more than one blue object to a scene and later responded that
922 *at least one* object must be blue. Thus, it seems that children's rules simply picked out
923 different features of the scenes than adults. An interesting question is whether, in the cases
924 where a child's guess is logically inconsistent with some of their learning data, this is
925 because their representation itself is imprecise, or because their verbal description
926 imprecisely describes their representation. Another possibility could be that adults are
927 better introspectors than children, better able to "read out" the structure of their own
928 representations (Morris, 2021). While these are intriguing possibilities our current
929 experiment cannot fully resolve these explanations.

930 **Implications for theories of concepts**

931 As we noted at the outset, psychological theories of concepts have oscillated between
932 symbolic accounts — that seek to explain conceptual productivity and creativity — and
933 similarity accounts — that seek to explain how concepts drive probabilistic generalization.
934 The constructivist framework is based in the symbolic camp, however it inherits many of
935 the advantages of similarity accounts by maintaining a relationship with probabilistic

936 inference embodied by the stochastic mechanisms of generation and search. Thus, we see
 937 our findings as support for recent claims that higher level cognition utilizes some form of
 938 stochastic generative sampling to approximate rational inference (Bramley, Dayan, et al.,
 939 2017; A. Sanborn et al., 2021; Zhu, Sanborn, & Chater, 2020) and that this might also
 940 explain aspects of human cultural and technological development that take place over
 941 populations and multiple generations (Krafft, Shmueli, Griffiths, Tenenbaum, et al., 2021).

942 As it stands, neither our PCFG or IDG are plausible process models of concept
 943 formation. The PCFG is a framework for normative top-down inference in the limit of
 944 infinite sampling, and IDG is a hybrid that is less sample inefficient as a brute force
 945 approach to inference in situations where a learner already has some evidence. A process
 946 account needs to explain how a learner searches the latent posterior in either framework
 947 with limited memory and computation. Following a number of recent research lines
 948 (Bramley, Mayrhofer, Gerstenberg, & Lagnado, 2017; Dasgupta, Schulz, & Gershman,
 949 2017; Fränken, Theodoropoulos, & Bramley, in revision; Ullman, Goodman, & Tenenbaum,
 950 2012), we see incremental adaptation of one or a few focal hypotheses in the light of
 951 evidence as a promising approach. A learner might use an observation to generate an
 952 initial idea akin to our IDG, but then explore permutations to this to account for the rest
 953 of the evidence. While hypothesis search models like RULEX (Nosofsky & Palmeri, 1998;
 954 Nosofsky et al., 1994) provide candidate heuristics for achieving such a search, their long
 955 run behavior lacks a clear relationship with computational-level rationality (Navarro,
 956 2005). However, if a learners' algorithmic adaptations approximate a valid approximation
 957 scheme, for instance accepting permutations with the Metropolis-Hastings probability
 958 $\max(1, \frac{P(r')}{P(r)})$ (Bramley, Dayan, et al., 2017; Dasgupta et al., 2016; Hastings, 1970; Thaker
 959 et al., 2017), we can explain why more probable hypotheses are discovered more often as
 960 well as explaining probability matching and order effects directly from our account of
 961 generation. Since the endpoint of an MCMC search approaches an independent posterior
 962 sample, we would expect a population of searchers to end up with a set of hypotheses that
 963 look like posterior samples. Moreover, since individual searchers have finite time to search,
 964 we would expect order effects and dependence in their ideas over time. To the extent that
 965 participants deviate from a probabilistically valid approximation scheme, for instance "hill
 966 climbing" or accepting only strictly better fitting ideas, we might also explain how they can
 967 get stuck in local optima and exhibit mal-adaptive order effects like garden paths (Gelpi,
 968 Prystawski, Lucas, & Buchsbaum, 2020).

969 Our modeling of generalizations revealed that there is no straightforward family
 970 resemblance between the features of rule-following training scenes (generated by the
 971 participant) and rule-following generalization scenes (as pre-selected for the experiment).

972 This resulted in our Similarity model performing at chance and uncorrelated with
973 participants while all our symbolic model variants received some support. While this is far
974 from an exhaustive test of sub-symbolic concept models, even a successful similarity-based
975 account of generalizations would only account for half of the behavior in this task. As well
976 as generalizing systematically, participants gave detailed natural language descriptions of
977 their ideas. The majority of these we could convert into logical statements (86%) that
978 predicted most generalizations (72%: children, 84%: adults) and were consistent with the
979 majority of their learning data (71%: children, 87%: adults). Any subsymbolic account of
980 concepts would essentially need to be paired with an explanation for how people generate
981 these verbal descriptions of their non-symbolic concepts that nonetheless reflect their use
982 (cf. Dennett, 1988). Arguably, this task is no easier than the one of generating a symbolic
983 hypothesis about the nature of the world in the first place. Thus we feel that our results
984 are more straightforwardly explained by our symbolic account whereby the logical
985 structure of the hypotheses participants describe is actually the causal mechanism driving
986 their generalizations rather than some form of computationally expensive but behaviorally
987 impotent retrospective confabulation (cf. Johansson et al., 2008).

988 Implications for theories of development

989 Our analyses revealed a variety of developmental differences. Children's guesses were
990 more complex than adults', and consequently we could capture them with a significantly
991 "flatter" generation process that inherently produced a wider diversity of hypotheses. This
992 is broadly normative: Having been exposed to less evidence, with less idea what conceptual
993 compositions and fragments will be useful in understanding their environment, we would
994 expect children's construction process to be less fine tuned. In other words, children are
995 justified in entertaining a wider set of ideas than adults. However, we noted at the start
996 that there are several potentially distinct algorithmic stories that could underpin this
997 diversity: (1) children might simply have hypothesis generation mechanism that embodies
998 a rationally flatter latent prior, (2) they might additionally explore more radically, over and
999 above differences in the relative credibility their latent prior actually attaches to different
1000 possibilities (Gopnik, 2020; Lucas et al., 2014; Wu et al., 2017) or (3) we also considered
1001 that children's generation mechanisms might be more dominated by "bottom-up"
1002 processes. We take our comparison of PCFG and IDG to speak against option 3. Adults'
1003 hypotheses were, as far as we could tell, at least as anchored to idiosyncratic patterns of
1004 their learning data as children's. However these data do not distinguish clearly between
1005 options (1) and (2). To do this, one would need to collect multiple judgments from the
1006 same participant within a learning problem, and somehow also elicit judgments about the

1007 relative quality of these guesses. If children's guesses shift within a problem in a way that
1008 is less sensitive to the relative subjective probabilities than adults, this would support the
1009 idea that children's hypothesis generation is more "high temperature" exploratory than
1010 adults' (Gopnik, 2020), over and above differences in the flatness of their latent prior.

1011 As well as producing more complex guesses, children also produced more elaborate
1012 scenes during active learning than adults. One possible characterization is that children's
1013 active scene construction also depended on a "flatter" generative prior, resulting in more
1014 diversity of exploration approaches. Indeed, differences in active exploration are the other
1015 side of the coin of the high temperature search idea (Friston et al., 2016; Gopnik, 2020;
1016 Klahr & Dunbar, 1988; E. Schulz, Klenske, Bramley, & Speekenbrink, 2017). On the other
1017 hand, the problem of generating informative tests is not quite the same as that of finding
1018 the right hypothesis. It is important to avoid redundancy and produce tests that target
1019 uncertainty and, in combination, serve to test a wide variety of salient hypotheses. In this
1020 sense, adults' testing behavior was more systematic, better reducing global measures of
1021 entropy and potentially reflecting a more top-down, or strategic, *control of variables*
1022 approach to gathering evidence and updating beliefs (Kuhn & Brannock, 1977). Within
1023 each trial, children's testing was more repetitive than adults', suggesting that they made
1024 slower progress in exploring the problem space, or were generally less able to keep track of
1025 what they had done, perhaps mixing more random actions in with systematic ones (Meder,
1026 Wu, Schulz, & Ruggeri, 2021).

1027 Children's guesses were also less consistent with their evidence than adults'. This
1028 might be because they were less able to extract common features across learning scenes
1029 (Ruggeri & Feufel, 2015; Ruggeri & Lombrozo, 2015). However, it could also be a
1030 consequence of a more generalized limitation in ability to generate, store and compare
1031 hypotheses. With a flatter prior and limited sampling, one has a lower chance generating a
1032 hypothesis that can explain all the evidence. Children also under-generalized, often
1033 selecting only 1 or 2 of the 8 test scenes (there was actually always 4) doing so even when
1034 their symbolic guesses predicted more should be selected. It could be that children found
1035 this part of the task overwhelming, perhaps tending to stop after identifying one or two
1036 hypothesis consistent scenes rather than evaluating all of them. On the face of it, this
1037 reflects Wu et al.'s (Wu et al., 2017) finding that children are weaker generalizers than
1038 adults.

1039 Curiously, children were more likely to refer to relational and positional properties
1040 in their guesses, while adults were by most likely to make guesses that pertained to the
1041 primary object features (color and size). This is an independently interesting finding. Since
1042 relational features are structurally more complex than primitive features, we might have

1043 predicted they would be more readily evoked by adults. It could be that children bought in
1044 more to the scientific reasoning cover story, treating mechanistic explanations, such as that
1045 objects must touch or be positioned in particular ways to produce stars, as credible
1046 (Gelman, 2004). Conversely, adults may have been more likely to expect Gricean
1047 considerations to apply, e.g. that experimenters would likely set simple rules using salient
1048 but abstract features like color over perceptually ambiguous properties like position
1049 (Szollosi & Newell, 2020). However, it could also be the case that there are deeper
1050 differences between the experiences of children and adults that render structural features
1051 more relevant to children and surface features more relevant to adults.

1052 **Conclusions**

1053 We analyzed an experiment combining rich qualitative and quantitative measures of
1054 children's and adults' inductive inference. We found a number of developmental differences
1055 and demonstrated that we can make sense of these through the computational
1056 constructivism lens. Our results add empirical support and theoretical detail to recent
1057 characterizations of children as more diverse thinkers and active learners than adults, and
1058 bring us closer to a computational understanding of human learning across the lifespan.

1059

References

- 1060 Allen, K. R., Smith, K. A., & Tenenbaum, J. B. (2020). Rapid trial-and-error learning
1061 with simulation supports flexible tool use and physical reasoning. *Proceedings of the
1062 National Academy of Sciences*, 117(47), 29302–29310.
- 1063 Bonawitz, E. B., Denison, S., Gopnik, A., & Griffiths, T. L. (2014). Win-stay, lose-sample:
1064 A simple sequential algorithm for approximating Bayesian inference. *Cognitive
1065 Psychology*, 74, 35–65.
- 1066 Box, G. E. (1976). Science and statistics. *Journal of the American Statistical Association*,
1067 71(356), 791–799.
- 1068 Bramley, N. R., Dayan, P., Griffiths, T. L., & Lagnado, D. A. (2017). Formalizing
1069 Neurath's ship: Approximate algorithms for online causal learning. *Psychological
1070 Review*, 124(3), 301–338.
- 1071 Bramley, N. R., Lagnado, D. A., & Speekenbrink, M. (2015). Conservative forgetful
1072 scholars: How people learn causal structure through interventions. *Journal of
1073 Experimental Psychology: Learning, Memory & Cognition*, 41(3), 708–731.
- 1074 Bramley, N. R., Mayrhofer, R., Gerstenberg, T., & Lagnado, D. A. (2017). Causal learning
1075 from interventions and dynamics in continuous time. In *Proceedings of the 39th*
1076 *Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive Science
1077 Society.
- 1078 Bramley, N. R., Rothe, A., Tenenbaum, J. B., Xu, F., & Gureckis, T. M. (2018).
1079 Grounding compositional hypothesis generation in specific instances. In *Proceedings
1080 of the 40th Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive
1081 Science Society.
- 1082 Brooks, P. J., & Braine, M. D. (1996). What do children know about the universal
1083 quantifiers all and each? *Cognition*, 60(3), 235–268.
- 1084 Bruner, J. S., Goodnow, J. J., & Austin, G. A. (1956). *A study of thinking*. Routledge.
- 1085 Bruner, J. S., Jolly, A., & Sylva, K. (1976). *Play: Its role in development and evolution*.
1086 Penguin.
- 1087 Buchanan, D. W., Tenenbaum, J. B., & Sobel, D. M. (2010). Edge replacement and
1088 nonindependence in causation. In *Proceedings of the 32nd Annual Meeting of the
1089 Cognitive Science Society* (pp. 919–924). Austin, TX: Cognitive Science Society.
- 1090 Carey, S. (1985). Are children fundamentally different kinds of thinkers and learners than
1091 adults. *Thinking and learning skills*, 2, 485–517.
- 1092 Carey, S. (2009). *The origin of concepts: Oxford series in cognitive development*. Oxford
1093 University Press, England.
- 1094 Chen, Z., & Klahr, D. (1999). All other things being equal: Acquisition and transfer of the

- 1095 control of variables strategy. *Child Development*, 70(5), 1098–1120.
- 1096 Clark, A. (2012). Whatever next? predictive brains, situated agents, and the future of
1097 cognitive science. *Behavioral Brain Sciences*, 1–86.
- 1098 Coenen, A., Nelson, J., & Gureckis, M., Todd. (2018). Asking the right questions about
1099 the psychology of human inquiry: Nine open challenges. *Psychonomic Bulletin &
1100 Review*, 26, 1548–1587.
- 1101 Coenen, A., Rehder, R., & Gureckis, T. M. (2015). Strategies to intervene on causal
1102 systems are adaptively selected. *Cognitive Psychology*, 79, 102–133.
- 1103 Coenen, A., Ruggeri, A., Bramley, N. R., & Gureckis, T. M. (2019). Testing one or
1104 multiple: How beliefs about sparsity affect causal experimentation. *Journal of
1105 Experimental Psychology: Learning, Memory, and Cognition*, 45(11), 1923.
- 1106 Cook, C., Goodman, N. D., & Schulz, L. E. (2011). Where science starts: Spontaneous
1107 experiments in preschoolers' exploratory play. *Cognition*, 120(3), 341–349.
- 1108 Dasgupta, I., Schulz, E., & Gershman, S. J. (2016). Where do hypotheses come from?
1109 *Center for Brains, Minds and Machines (preprint)*.
- 1110 Dasgupta, I., Schulz, E., & Gershman, S. J. (2017). Where do hypotheses come from?
1111 *Cognitive psychology*, 96, 1–25.
- 1112 Dennett, D. C. (1988). The intentional stance in theory and practice. In R. Byrne &
1113 A. Whiten (Eds.), *Machiavellian intelligence* (pp. 180–202). Oxford, UK: Oxford
1114 University Press.
- 1115 Dennett, D. C. (1991). *Consciousness explained*. London, UK: Penguin.
- 1116 Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., ... Tenenbaum, J. B.
1117 (2020). Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep
1118 bayesian program learning. *arXiv preprint arXiv:2006.08381*.
- 1119 Fedyk, M., & Xu, F. (2018). The epistemology of rational constructivism. *Review of
1120 Philosophy and Psychology*, 9(2), 343–362.
- 1121 Feldman, J. (2000). Minimization of Boolean complexity in human concept learning.
1122 *Nature*, 407(6804), 630.
- 1123 Fodor, J. A. (1975). *The language of thought* (Vol. 5). Harvard university press.
- 1124 Fränken, J.-P., Theodoropoulos, N. C., & Bramley, N. R. (in revision). Algorithms for
1125 adaptation in inductive inference.
- 1126 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., Pezzulo, G., et al. (2016). Active
1127 inference and learning. *Neuroscience & Biobehavioral Reviews*, 68, 862–879.
- 1128 Gelman, S. A. (2004). Psychological essentialism in children. *Trends in cognitive sciences*,
1129 8(9), 404–409.
- 1130 Gelpi, R., Prystawski, B., Lucas, C. G., & Buchsbaum, D. (2020). Incremental hypothesis

- revision in causal reasoning across development.
- Gershman, S. J., & Niv, Y. (2010). Learning latent structure: carving nature at its joints. *Current Opinion in Neurobiology*, 20(2), 251–256.
- Gettys, C. F., & Fisher, S. D. (1979). Hypothesis plausibility and hypothesis generation. *Organizational behavior and human performance*, 24(1), 93–110.
- Ginsburg, S. (1966). *The mathematical theory of context free languages*. McGraw-Hill Book Company.
- Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational analysis of rule-based concept learning. *Cognitive Science*, 32(1), 108–154.
- Goodman, N. D., Ullman, T. D., & Tenenbaum, J. B. (2011). Learning a theory of causality. *Psychological Review*, 118(1), 110–9.
- Gopnik, A. (1996). The scientist as child. *Philosophy of science*, 63(4), 485–514.
- Gopnik, A. (2020). Childhood as a solution to explore-exploit tensions. *Philosophical Transactions of the Royal Society B*, 375(1803), 20190502.
- Gopnik, A., Glymour, C., Sobel, D., Schulz, L. E., Kushnir, T., & Danks, D. (2004). A theory of causal learning in children: Causal maps and Bayes nets. *Psychological Review*, 111, 1–31.
- Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in Cognitive Science*, 7, 217–229.
- Griffiths, T. L., & Tenenbaum, J. B. (2009). Theory-based causal induction. *Psychological Review*, 116, 661–716.
- Gureckis, T. M., & Markant, D. B. (2012, September). Self-Directed Learning: A Cognitive and Computational Perspective. *Perspectives on Psychological Science*, 7(5), 464–481.
- Hastings, W. K. (1970). Monte carlo sampling methods using markov chains and their applications.
- Heath, C. (2004). *Zendo—Design History*. Retrieved from <http://www.koryheath.com/zendo/design-history/>
- Howson, C., & Urbach, P. (2006). *Scientific reasoning: the Bayesian approach*. Open Court Publishing.
- Inhelder, B., & Piaget, J. (1958). *The growth of logical thinking from childhood to adolescence: An essay on the construction of formal operational structures* (Vol. 22). Psychology Press.
- Johansson, P., Hall, L., & Sikström, S. (2008). From change blindness to choice blindness. *Psychologica*, 51(2), 142–155.

- 1167 Johnson, R. B., Onwuegbuzie, A. J., & Turner, L. A. (2007). Toward a definition of mixed
1168 methods research. *Journal of mixed methods research*, 1(2), 112–133.
- 1169 Johnson-Laird, P. N. (1983). *Mental models: Towards a cognitive science of language,
1170 inference, and consciousness*. Cambridge: Cambridge University Press.
- 1171 Jones, A., Bramley, N. R., Gureckis, T. M., & Ruggeri, A. (in revision). Changing many
1172 things at once sometimes makes for a good experiment, and children know that.
- 1173 Jones, M., & Love, B. C. (2011). Bayesian Fundamentalism or Enlightenment? On the
1174 explanatory status and theoretical contributions of Bayesian models of cognition. *The
1175 Behavioral and Brain Sciences*, 34(4), 169–88.
- 1176 Kemp, C., Goodman, N. D., & Tenenbaum, J. B. (2010). Learning to learn causal models.
1177 *Cognitive Science*, 34(7), 1185–243.
- 1178 Kemp, C., & Tenenbaum, J. B. (2009). Structured statistical models of inductive
1179 reasoning. *Psychological Review*, 116(1), 20.
- 1180 Klahr, D., & Dunbar, K. (1988). Dual space search during scientific reasoning. *Cognitive
1181 Science*, 12(1), 1–48.
- 1182 Klahr, D., Fay, A. L., & Dunbar, K. (1993). Heuristics for scientific experimentation: A
1183 developmental study. *Cognitive Psychology*, 25(1), 111–146.
- 1184 Klahr, D., Zimmerman, C., & Jirout, J. (2011). Educational interventions to advance
1185 children's scientific thinking. *Science*, 333(6045), 971–975.
- 1186 Klayman, J., & Ha, Y.-w. (1989). Hypothesis testing in rule discovery: Strategy, structure,
1187 and content. *Journal of Experimental Psychology: Learning, Memory & Cognition*,
1188 15(4), 596.
- 1189 Komatsu, L. K. (1992). Recent views of conceptual structure. *Psychological bulletin*,
1190 112(3), 500.
- 1191 Krafft, P. M., Shmueli, E., Griffiths, T. L., Tenenbaum, J. B., et al. (2021). Bayesian
1192 collective learning emerges from heuristic social learning. *Cognition*, 212, 104469.
- 1193 Krippendorff, K. (2012). *Content analysis: An introduction to its methodology*. Sage.
- 1194 Kruschke, J. K. (1992). Alcove: an exemplar-based connectionist model of category
1195 learning. *Psychological Review*, 99(1), 22.
- 1196 Kuhn, D., & Brannock, J. (1977). Development of the isolation of variables scheme in
1197 experimental and “natural experiment” contexts. *Developmental psychology*, 13(1),
1198 9.
- 1199 Lagnado, D. A., & Sloman, S. A. (2006). Time as a guide to cause. *Journal of
1200 Experimental Psychology: Learning, Memory & Cognition*, 32(3), 451–60.
- 1201 Lai, L., & Gershman, S. J. (2021). Policy compression: an information bottleneck in action
1202 selection.

- 1203 Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building
1204 machines that learn and think like people. *Behavioral and Brain Sciences*, 40.
- 1205 Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and
1206 reversals. In *Soviet physics doklady* (Vol. 10, pp. 707–710).
- 1207 Lewis, O., Perez, S., & Tenenbaum, J. (2014). Error-driven stochastic search for theories
1208 and concepts. In *Proceedings of the annual meeting of the cognitive science society*
1209 (Vol. 36).
- 1210 Lieder, F., Griffiths, T., Huys, Q. J., & Goodman, N. D. (2017). The anchoring bias
1211 reflects rational use of cognitive resources.
- 1212 Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). Sustain: a network model of category
1213 learning. *Psychological Review*, 111(2), 309.
- 1214 Lucas, C. G., Bridgers, S., Griffiths, T. L., & Gopnik, A. (2014). When children are better
1215 (or at least more open-minded) learners than adults: Developmental differences in
1216 learning the forms of causal relationships. *Cognition*, 131(2), 284–299.
- 1217 Lucas, C. G., & Griffiths, T. L. (2010). Learning the form of causal relationships using
1218 hierarchical bayesian models. *Cognitive Science*, 34(1), 113–147.
- 1219 Luce, D. R. (1959). *Individual choice behavior*. New York: Wiley.
- 1220 Markant, D. B., & Gureckis, T. M. (2014). Is it better to select or to receive? learning via
1221 active and passive hypothesis testing. *Journal of Experimental Psychology: General*,
1222 143(1), 94.
- 1223 Marr, D. (1982). *Vision*. New York: Freeman & Co.
- 1224 McCormack, T., Bramley, N. R., Frosch, C., Patrick, F., & Lagnado, D. A. (2016).
1225 Children's use of interventions to learn causal structure. *Journal of Experimental
1226 Child Psychology*, 141, 1–22.
- 1227 Meder, B., Wu, C. M., Schulz, E., & Ruggeri, A. (2021). Development of directed and
1228 random exploration in children. *Developmental Science*, 24(4), e13095.
- 1229 Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning.
1230 *Psychological Review*, 85(3), 207.
- 1231 Meng, Y., Bramley, N., & Xu, F. (2018). Children's causal interventions combine
1232 discrimination and confirmation. In *Proceedings of the 40th annual conference of the
1233 cognitive science society*.
- 1234 Michalski, R. S. (1969). On the quasi-minimal solution of the general covering problem. In
1235 (Vol. A3, pp. 125–128).
- 1236 Morris, A. (2021). Invisible gorillas in the mind: Internal inattentional blindness and the
1237 prospect of introspection training.
- 1238 Navarro, D. J. (2005). Analyzing the rulex model of category learning. *Journal of*

- 1239 *Mathematical Psychology*, 49(4), 259–275.
- 1240 Nelson, J. D. (2005). Finding useful questions: on Bayesian diagnosticity, probability,
1241 impact, and information gain. *Psychological Review*, 112(4), 979–99.
- 1242 Nosofsky, R. M., & Palmeri, T. J. (1998). A rule-plus-exception model for classifying
1243 objects in continuous-dimension spaces. *Psychonomic Bulletin & Review*, 5(3),
1244 345–369.
- 1245 Nosofsky, R. M., Palmeri, T. J., & McKinley, S. C. (1994). Rule-plus-exception model of
1246 classification learning. *Psychological review*, 101(1), 53.
- 1247 Nye, M. I., Solar-Lezama, A., Tenenbaum, J. B., & Lake, B. M. (2020). Learning
1248 compositional rules via neural program synthesis. *arXiv preprint arXiv:2003.05562*.
- 1249 Oaksford, M., & Chater, N. (2007). *Bayesian rationality: the probabilistic approach to*
1250 *human reasoning*. Oxford: Oxford University Press.
- 1251 Perfors, A., Tenenbaum, J. B., Griffiths, T. L., & Xu, F. (2011). A tutorial introduction to
1252 bayesian models of cognitive development. *Cognition*, 120(3), 302–321.
- 1253 Piaget, J. (2013). *The construction of reality in the child* (Vol. 82). Routledge.
- 1254 Piaget, J., & Valsiner, J. (1930). *The child's conception of physical causality*. Transaction
1255 Pub.
- 1256 Piantadosi, S. T. (2021). The computational origin of representation. *Minds and*
1257 *Machines*, 31(1), 1–58.
- 1258 Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2012). Bootstrapping in a
1259 language of thought: A formal model of numerical concept learning. *Cognition*,
1260 123(2), 199–217.
- 1261 Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2016). The logical primitives of
1262 thought: Empirical foundations for compositional cognitive models. *Psychological*
1263 *Review*, 123(4), 392.
- 1264 Posner, M. I., & Keele, S. W. (1968). On the genesis of abstract ideas. *Journal of*
1265 *experimental psychology*, 77(3p1), 353.
- 1266 Quine, W. v. O. (1969). *Word and object*. MIT press.
- 1267 Rothe, A., Lake, B. M., & Gureckis, T. M. (2017). Question asking as program generation.
1268 *arXiv preprint arXiv:1711.06351*.
- 1269 Ruggeri, A., & Feufel, M. (2015). How basic-level objects facilitate question-asking in a
1270 categorization task. *Frontiers in psychology*, 6, 918.
- 1271 Ruggeri, A., & Lombrozo, T. (2014). Learning by asking: How children ask questions to
1272 achieve efficient search. In *36th annual meeting of the cognitive science society* (pp.
1273 1335–1340). Austin, TX: Cognitive Science Society.
- 1274 Ruggeri, A., & Lombrozo, T. (2015). Children adapt their questions to achieve efficient

- 1275 search. *Cognition*, 143, 203–216.
- 1276 Ruggeri, A., Lombrozo, T., Griffiths, T. L., & Xu, F. (2016). Sources of developmental
1277 change in the efficiency of information search. *Developmental psychology*, 52(12),
1278 2159.
- 1279 Ruis, L., Andreas, J., Baroni, M., Bouchacourt, D., & Lake, B. M. (2020). A benchmark
1280 for systematic generalization in grounded language understanding. *arXiv preprint*
1281 *arXiv:2003.05161*.
- 1282 Rule, J. S., Schulz, E., Piantadosi, S. T., & Tenenbaum, J. B. (2018). Learning list
1283 concepts through program induction. *BioRxiv*, 321505.
- 1284 Rule, J. S., Tenenbaum, J. B., & Piantadosi, S. T. (2020). The child as hacker. *Trends in*
1285 *Cognitive Sciences*.
- 1286 Sanborn, A., Zhu, J., Spicer, J., Sundh, J., León-Villagrá, P., & Chater, N. (2021).
1287 Sampling as the human approximation to probabilistic inference.
- 1288 Sanborn, A. N., & Chater, N. (2016). Bayesian brains without probabilities. *Trends in*
1289 *Cognitive Sciences*.
- 1290 Schulz, E., Klenske, E. D., Bramley, N. R., & Speekenbrink, M. (2017). Strategic
1291 exploration in human adaptive control. *bioRxiv*, 110486.
- 1292 Schulz, L. E., Goodman, N. D., Tenenbaum, J. B., & Jenkins, A. C. (2008). Going beyond
1293 the evidence: Abstract laws and preschoolers' responses to anomalous data.
1294 *Cognition*, 109(2), 211–223.
- 1295 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2),
1296 461–464.
- 1297 Shanks, D. R., Tunney, R. J., & McCarthy, J. D. (2002). A re-examination of probability
1298 matching and rational choice. *Journal of Behavioral Decision Making*, 15(3),
1299 233–250.
- 1300 Shannon, C. E. (1951). Prediction and entropy of printed english. *The Bell System*
1301 *Technical Journal*, 30, 50–64.
- 1302 Shepard, R. N. (1987). Toward a universal law of generalization for psychological science.
1303 *Science*, 237(4820), 1317–1323.
- 1304 Shepard, R. N., & Chang, J.-J. (1963). Stimulus generalization in the learning of
1305 classifications. *Journal of Experimental Psychology*, 65(1), 94.
- 1306 Sim, Z. L., & Xu, F. (2017). Learning higher-order generalizations through free play:
1307 Evidence from 2-and 3-year-old children. *Developmental psychology*, 53(4), 642.
- 1308 Sobel, D. M., & Kushnir, T. (2006). The importance of decision making in causal learning
1309 from interventions. *Memory & Cognition*, 34(2), 411–419.
- 1310 Stewart, N., Chater, N., & Brown, G. D. A. (2006). Decision by sampling. *Cognitive*

- 1311 *Psychology*, 53(1), 1–26.
- 1312 Steyvers, M., Tenenbaum, J. B., Wagenmakers, E., & Blum, B. (2003). Inferring causal
1313 networks from observations and interventions. *Cognitive Science*, 27, 453–489.
- 1314 Szollosi, A., & Newell, B. R. (2020). People as intuitive scientists: Reconsidering statistical
1315 explanations of decision making. *Trends in Cognitive Sciences*.
- 1316 Tenenbaum, J. B. (1999). *A bayesian framework for concept learning* (Unpublished
1317 doctoral dissertation). Massachusetts Institute of Technology.
- 1318 Thaker, P., Tenenbaum, J. B., & Gershman, S. J. (2017). Online learning of symbolic
1319 concepts. *Journal of Mathematical Psychology*, 77, 10–20.
- 1320 Turing, A. M. (1937). On computable numbers, with an application to the
1321 entscheidungsproblem. *Proceedings of the London mathematical society*, 2(1),
1322 230–265.
- 1323 Tversky, A. (1977). Features of similarity. *Psychological Review*, 84(4), 327.
- 1324 Ullman, T. D., Goodman, N. D., & Tenenbaum, J. B. (2012). Theory learning as stochastic
1325 search in the language of thought. *Cognitive Development*, 27(4), 455–480.
- 1326 Valkov, L., Chaudhari, D., Srivastava, A., Sutton, C., & Chaudhuri, S. (2018). Houdini:
1327 Lifelong learning as program synthesis. In *Advances in Neural Information
1328 Processing Systems* (pp. 8687–8698).
- 1329 Van Laarhoven, P. J., & Aarts, E. H. (1987). Simulated annealing. In *Simulated annealing:
1330 Theory and applications* (pp. 7–15). Springer.
- 1331 Van Rooij, I., Blokpoel, M., Kwisthout, J., & Wareham, T. (2019). *Cognition and
1332 intractability: A guide to classical and parameterized complexity analysis*. Cambridge
1333 University Press.
- 1334 von Humboldt, W. (1863/1988). *On language*. New York: Cambridge University Press.
- 1335 Vul, E., Goodman, N. D., Griffiths, T. L., & Tenenbaum, J. B. (2009). One and done?
1336 optimal decisions from very few samples. In *Proceedings of the 31st Annual Meeting
1337 of the Cognitive Science Society* (Vol. 1, pp. 66–72). Austin, TX: Cognitive Science
1338 Society.
- 1339 Wason, P. C. (1960). On the failure to eliminate hypotheses in a conceptual task.
1340 *Quarterly journal of experimental psychology*, 12(3), 129–140.
- 1341 Wason, P. C. (1968). Reasoning about a rule. *The Quarterly Journal of Experimental
1342 Psychology*, 20(3), 273–281.
- 1343 Wu, C. M., Schulz, E., Speekenbrink, M., Nelson, J. D., & Meder, B. (2017). Exploration
1344 and generalization in vast spaces. *bioRxiv*. doi: 10.1101/171371
- 1345 Xu, F. (2019). Towards a rational constructivist theory of cognitive development.
1346 *Psychological Review*, 126(6), 841.

- 1347 Zhu, J.-Q., Sanborn, A. N., & Chater, N. (2020). The bayesian sampler: Generic bayesian
1348 inference causes incoherence in human probability judgments. *Psychological review*,
1349 127(5), 719.

1350

Appendix A: Models

1351 Generating context free (PCFG) model predictions

1352 We created a grammar (specifically a *probabilistic context free grammar* or PCFG;
 1353 Ginsburg, 1966) that can be used to produce any rule that can be expressed with
 1354 first-order logic and lambda abstraction referring to the features participants referred to in
 1355 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...	$N_{\{<, >, =\}}(\lambda x_i : .., N, \mathcal{X})$
Feature f of x_i has value {larger, smaller, (or) equal} to v	$\{<, >, \leq, \geq, =\}(x_i, v, f)$
Feature f of x_i is {larger, smaller, (or) equal} to feature f of x_j	$\{<, >, \leq, \geq, =\}(x_i, x_j, f)$
Relation r between x_i and x_j holds	$\Gamma(x_i, x_j, r)$
Booleans {and,or,not}	$\{\wedge, \vee, \neq\}(x)$
Object feature	Levels
Color	{red, green, blue}
Size	{1:small, 2:medium, 3:large}
x -position	(0,8)
y -position	(0,8)
Orientation	{Upright, left hand side, right hand side, strange}
Grounded	true if touching the ground
Pairwise feature	Condition
Contact	true if x_1 touches x_2
Stacked	true if x_1 is above and touching x_2 and x_2 is grounded
Pointing	true if x_1 is orientated {left/right} and x_2 is to x_1 's {left/right}
Inside	true if x_1 is smaller than x_2 + has same x and y position (± 0.3), false

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

1356

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

1357 to set up a set of string-rewrite rules (Goodman et al., 2008). Thus, each hypothesis begins
 1358 life as a string containing a single *non-terminal symbol* (here, S) that is replaced using

1359 rewrite rules, or *productions*. These productions are repeatedly applied to the string,
 1360 replacing non-terminal symbols with a mixture of other non-terminal symbols and terminal
 1361 fragments of first order logic, until no non-terminal symbols remain. The productions are
 1362 so designed that the resulting string is guaranteed to be a valid grammatical expression
 1363 and all grammatical expressions have a nonzero chance of being produced. In addition, by
 1364 having the productions tie the expression to bound variables and truth statements, our
 1365 PCFG serves as an automatic concept generator. Table A-2 details the PCFG we used in
 1366 the paper.

1367 We use capital letters as non-terminal symbols and each rewrite is sampled from the
 1368 available productions for a given symbol.⁶ Because some of the productions involve
 1369 branching (e.g., $B \rightarrow H(B, B)$), the resultant string can become arbitrarily long and
 1370 complex, involving multiple boolean functions and complex relationships between bound
 1371 variables.

1372 We include a variant that samples uniformly from the set of possible replacements
 1373 in each case, but we also reverse engineer a set of productions that produce exactly the
 1374 statistics of the empirical samples, as described in the main text.

1375 We used the process described in A-2 to produce a sample of 10,000 with a uniform
 1376 generation prior and an additional 10,000 for each participant with a “held out”
 1377 age-consistent prior based on the rule guesses of other participants in the requisite
 1378 agegroup. For the flipped prior analyses, we used the sample generated for the
 1379 chronologically first participant from the other agegroup.

1380 Generating instance driven (IDG) model predictions

1381 We used the algorithm proposed in Bramley et al. (2018) to produce a sample of
 1382 10,000 “grounded hypotheses” for each participant and trial, splitting these evenly across
 1383 the 8 learning scenes that participant produced and tested. For each, we generated two
 1384 sets: One using a uniform construction weights, and one with an age-appropriate “held
 1385 out” set of weights based on the rule guesses of other participants in the requisite agegroup.
 1386 For the flipped prior analyses, we used the weights from the chronologically first participant
 1387 from the other agegroup to generate samples inspired by the current participants’ evidence.

1388 To generate hypotheses as candidates for the hidden rule, the model uses the
 1389 following procedure with probabilities either set to uniform or drawn from the PCFG-fitted

⁶ 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 (numeric only)	$D1 \rightarrow$	value,	feature	
	$D2 \rightarrow$	value,	feature	
Feature (numeric only)	$E1 \rightarrow$	feature		
	$E2 \rightarrow$	feature		
(relational)	$E3 \rightarrow$	feature		
Boolean	$J \rightarrow$	\wedge	\vee	...
Inequality	$I \rightarrow$	\leq	\geq	$>$
		$<$		
Number	$K \rightarrow$	$n \in \{1, 2, 3, 4, 5, 6\}$		

Note: Context-sensitive aspects of the grammar: ^aBound variable(s) sampled uniformly without replacement from set; expressions requiring multiple variables censored if only one.

1390 productions for adults or for children (Figure 3g&h) and denoted with square brackets:

1391 1. **Observe.** either:

- 1392 (a) With probability $[A \rightarrow B]$: Sample a cone from the observation, then sample
 1393 one of its features f with probability $[G \rightarrow f]$ — e.g., $\{\#1\}$:⁷ “medium, size” or
 1394 $\{\#3\}$: “red, color”.
- 1395 (b) With probability $[A \rightarrow \text{Start}]$: Sample two cones uniformly without replacement
 1396 from the observation, and sample any shared or pairwise feature — e.g.,
 1397 $\{\#1, \#2\}$: “size”, or “contact”

1398 2. **Functionize.** Bind a variable for each sampled cone in Step 1 and sample a true
 1399 (in)equality statement relating the variable(s) and feature:

- 1400 (a) For a statement involving an unordered feature there is only one possibility —
 1401 e.g., $\{\#3\}$: “ $= (x_1, \text{red}, \text{color})$ ”, or for $\{\#1, \#2\}$: “ $= (x_1, x_2, \text{color})$ ”
- 1402 (b) For a single cone and an ordered feature, this could also be a nonstrict
 1403 inequality (\geq or \leq). We assume a learner only samples an inequality if it
 1404 expands the number of cones picked out from the scene relative to an equality

⁷ Numbers prepended with # refer to the labels on the cones in the example observation in Figure 2b.

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

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

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

- (a) **Conjunction.** A cone is sampled from the subset picked out by the statement thus far and one of its features sampled with probability $[G \rightarrow f]$ — e.g., $\{\#1\} \wedge (=x_1, \text{green}, \text{color}), \geq(x_1, \text{medium}, \text{size})$). Again, inequalities are sampleable only if they increase the true set size relative to equality — e.g., “ $\wedge(\leq(x_1, 3, \text{xposition}), \geq(x_1, \text{medium}, \text{size}))$ ”, which picks out more objects than “ $\wedge(=x_1, 3, \text{xposition}), \geq(x_1, \text{medium}, \text{size})$ ”.
- (b) **Disjunction.** An additional feature-value pair is selected uniformly from *either* unselected values of the current feature, *or* from a different feature — e.g., $\vee(=x_1, \text{color}, \text{red}), =x_1, \text{color}, \text{blue})$ or $\vee(=x_1, \text{color}, \text{blue}), \geq(x_1, \text{size}, 2)$). This step is skipped if the statement is already true of all the cones in the scene.⁸

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

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

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

⁸ We rounded positional features to one decimal place in evaluating rules to allow for perceptual uncertainty.

1435 $\forall(\lambda x_i : A, \mathcal{X})$ is censored and the probabilities re-normalized. K is set to match
 1436 number of cones for which the statement is true.

- 1437 (b) Statements involving two bound variables in lambda calculus have two nested
 1438 quantifier statements each selected as in (a). The inner statement quantifying x_2
 1439 is selected first based on truth value of the expression while taking x_1 to refer to
 1440 the cone observed in ‘1.’. The truth of the selected inner quantified statement is
 1441 then assessed for all cones to select the outer quantifier — e.g., $\{\#3, \#4\}$
 1442 “ $\wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size}))$ ” might become
 1443 “ $\forall(\lambda x_1 : \exists(\lambda x_2 : \wedge(= (x_2, \text{green}, \text{color}), \leq (x_1, x_2, \text{size})), \mathcal{X}), \mathcal{X})$ ”. The inner
 1444 quantifier \exists is selected (three of the four cones are green $\{\#1, \#2, \#4\}$), and
 1445 the outer quantifier \forall is selected (all cones are less than or equal in size to a
 1446 green cone).

1447 Note that a procedure like the one laid out above is, in principle, capable of
 1448 generating any rule generated by the PCFG in Figure 3g&3h, but will only do so when
 1449 exposed to an observation that is actually consistent with that rule, and will do so more
 1450 often when the observation is inconsistent with as many other rules as possible (i.e., a
 1451 minimal positive example). Step 4. allows that non-rule following scenes can be used to
 1452 inspire rules involving a negation, for instance that “something is not upright” – which is
 1453 semantically equivalent to saying that “nothing is upright”. Basing hypotheses on instances
 1454 may improve the quality of the effective sample of hypotheses that the learner generates.

1455 One way to think of the IDG procedure is as a partial inversion of a PCFG. As
 1456 illustrated by the blue text in the examples in Figure 2b in the main text. While the
 1457 PCFG starts at the outside and works inward, the IDG starts from the central content and
 1458 works outward out to a quantified statement, ensuring at each step that this final
 1459 statement is true of the scene.

1460 We note that it is possible, in principle, to calculate a lower bound on the prior
 1461 probability for the PCFG or IDG generating a hypothesis that a participant reported, even
 1462 if it does not occur in our sample. This can be achieved by reverse engineering the
 1463 production steps that would be needed to produce the precise encoded syntax. This is a
 1464 lower bound because it does not count semantically equivalent “phrasings” of the
 1465 hypothesis that e.g. mention features in different orders or use logically equivalent
 1466 combinations of booleans. We found that complex expressions tend to have a large number
 1467 of “phrasings”. In our sample-based approximation we implicitly treat semantically
 1468 equivalent expressions as constituting the same hypothesis but note that determining
 1469 semantic equivalence is an nontrivial aspect of constructivist inference that we do not fully

¹⁴⁷⁰ address here.

¹⁴⁷¹ **Full generalization model fits**

¹⁴⁷² As described in main text, we fit 18 model variants to participant’s data. All models
¹⁴⁷³ have between 0 and 2 parameters. For each model, we fit the parameter(s) by maximizing
¹⁴⁷⁴ the model’s likelihood of producing the participant data, using R’s `optim` function. We
¹⁴⁷⁵ compare models using the Bayesian Information Criterion (Schwarz, 1978) to accommodate
¹⁴⁷⁶ their different numbers of fitted parameters.⁹ Full results are in Table A-3.

¹⁴⁷⁷ **Appendix B: Free response coding**

¹⁴⁷⁸ To analyze the free responses, we first had two coders go through all responses and
¹⁴⁷⁹ categorize them as either:

- ¹⁴⁸⁰ 1. Correct: The subject gives exactly the correct rule or something logically equivalent
- ¹⁴⁸¹ 2. Overcomplicated: The subject gives a rule that over-specifies the criteria needed to
¹⁴⁸² produce stars relative to the ground truth. This means the rule they give is logically
¹⁴⁸³ sufficient but not necessary. For example, stipulating that “there must be a small
¹⁴⁸⁴ red” is overcomplicated if the true rule is “there must be a red” because a scene could
¹⁴⁸⁵ contain a medium or large red and emit stars.
- ¹⁴⁸⁶ 3. Overliberal: The opposite of overcomplicated. The subject gives a rule that
¹⁴⁸⁷ under-specifies what must happen for the scene to produce stars. For example,
¹⁴⁸⁸ stipulating that “there must be a blue” if the true rule is that “exactly one is blue”.
¹⁴⁸⁹ This is logically necessary but not sufficient because a scene could contain blue
¹⁴⁹⁰ objects but not produce stars because there is not exactly one of them.
- ¹⁴⁹¹ 4. Different: The subject gives a rule that is intelligible but different from the ground
¹⁴⁹² truth in that it is neither necessary or sufficient for determining whether a scene will
¹⁴⁹³ produce stars.
- ¹⁴⁹⁴ 5. Vague or multiple. Nuisance category.

⁹ On one perspective, our derivation of the child-like and adult-like productions constitutes fitting an additional 39 parameters ($m - 1$ for each production step), so evoking an additional BIC parameter penalty of $39 \times \log(3940) = 323$ for PCFG over PCFG Uniform and similarly for the IDG. If we were to apply this penalty, the uniform weighted variants would be clearly preferred under the BIC criterion at the aggregate level. It is less clear how to apply this penalty at the individual level 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	λ	τ	N	Accuracy
1. Baseline	children	-1319.75	2639.50			7	50%
2. Bias	children	-1218.96	2445.47	0.32		16	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%
18. Symbolic Guess + Bias	children	-1067.18	2149.47	0.26	0.80	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	32	70%
18. Symbolic Guess + Bias	adults	-880.59	1776.38	0.08	0.50	4	

NB: Accuracy column shows performance of the requisite model across 100 simulated runs through the task using participants' active learning data with τ 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.

1495 6. No rule. The subject says they cannot think of a rule.

1496 We were able to encode 205/238 (86%) of the children's responses and (219/250)
 1497 87% for adults as correct, overcomplicated, overliberal or different. Table A-4 shows the
 1498 complete confusion matrix. The two coders agreed 85% of the time, resulting in a Cohen's
 1499 Kappa of .77 indicating a good level of agreement (Krippendorff, 2012).

1500 We then had one coder familiar with the grammar go through each free response
 1501 that was not assigned vague or no rule, and encode it as a function in our grammar. The
 1502 second coder then blind spot checked 15% of these rules (64) and agreed in 95% of cases
 1503 61/64. The 6 cases of disagreement were discussed and resolved. In 5/6 cases, this was in
 1504 favor of the primary coder. The full set of free text responses along with the requisite

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

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

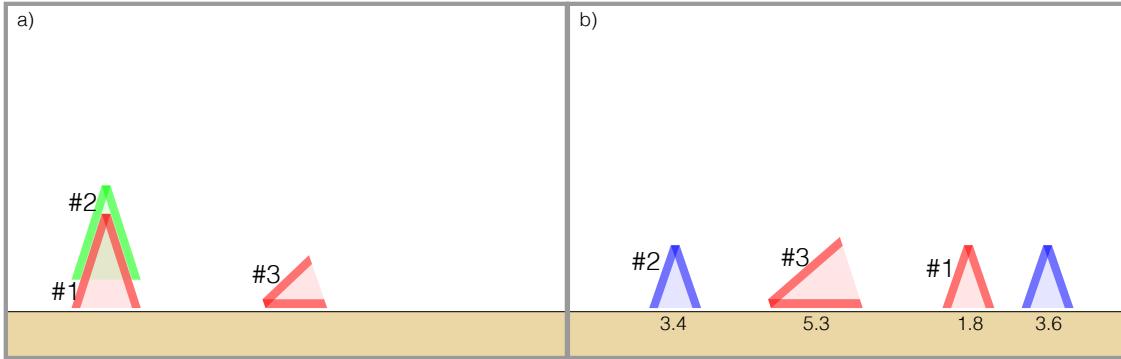
1505 classification, encoded rules are available in the [Online Repository](#).

1506 **Appendix C: Scene similarity measurement**

1507 To establish the overall similarity between two scenes, we need to map the objects
 1508 in a given scene to the objects in another scene (for example between the scenes in
 1509 FigureA-1 a and b) and establish a reasonable cost for the differences between objects
 1510 across dimensions. We also need a procedure for cases where there are objects in one scene
 1511 that have no analogue in the other. We approach the calculation of similarity via the
 1512 principle of minimum edit distance (Levenshtein, 1966). This means summing up the
 1513 elementary operations required to convert scene (a) into scene (b) or visa versa. We assume
 1514 objects can be adjusted in one dimension at a time (i.e. moving them on the x axis,
 1515 rotating them, or changing their color, and so on).

1516 Before focusing on how to map the objects between the scenes we must decide how
 1517 to measure the adjustment distance for a particular object in scene a to its supposed
 1518 analogue in scene b. As a simple way to combine the edit costs across dimensions we first
 1519 Z-score each dimension, such that the average distance between any two values across all
 1520 objects and all scenes and dimensions is 1. We then take the L1-norm (or city block
 1521 distance) as the cost for converting an object in scene (a) to an object in scene (b), or visa
 1522 versa. Note this is sensitive the size of the adjustment, penalizing larger changes in
 1523 position, orientation or size more severely than smaller changes, while changes in color are
 1524 all considered equally large since color is taken as categorical. Note also that for
 1525 orientation differences we also always assume the shortest distance around the circle.

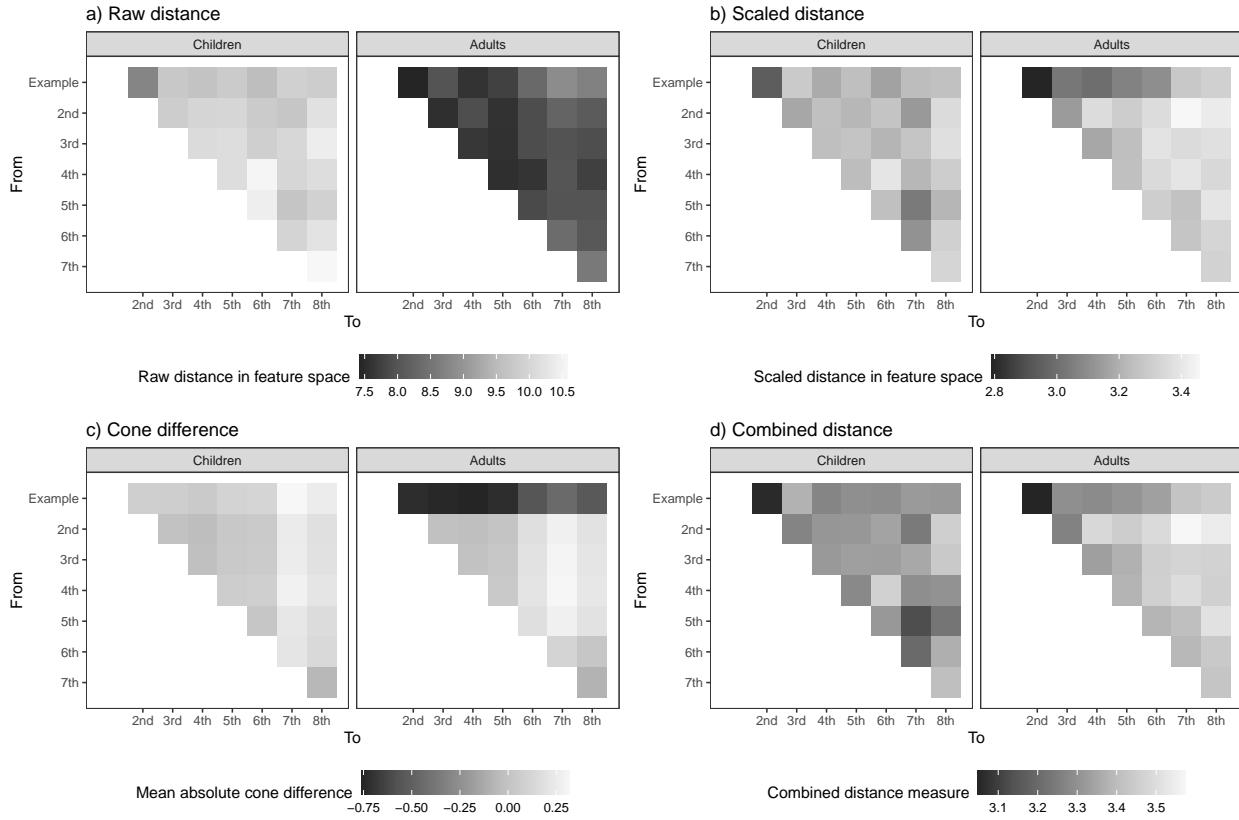
1526 If scene (a) has an object that does not exist in scene (b) we assume a default
 1527 adjustment penalty equal to the average divergence between two objects across all
 1528 comparisons (3.57 in the current dataset). We do the same for any object that exists in (a)
 1529 but not (b).

**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). Intuitively scene a) is more similar to scene b) than to scene c) and this is reflected in the similarity scores.

1530 Calculating the overall similarity between two scenes involves solving a mapping
 1531 problem of identifying which objects in scene (a) are “the same” as those in scene (b). We
 1532 resolve this “charitably”, by searching exhaustively for the mapping of objects in scene (a)
 1533 to scene (b) that minimizes the total edit distance. Having selected this mapping, and
 1534 computed the final edit distance including any costs for additional or removed objects, we
 1535 divide by the number shared cones, so as to avoid the dissimilarities increasing with the
 1536 number of objects involved.

1537 Figure A-2 computes the inter-scene similarity components that go into Figure 6c in
 1538 the main text. Summing up the edit distances across all objects, children’s scenes seem
 1539 much more diverse than adults (Figure A-2a). However this is primarily due to their
 1540 containing a greater average number of objects. Scaling the edit distance by the number of
 1541 objects in the target scene gives a more balanced perspective (Figure A-2b) but does not
 1542 account for the fact that the compared scene may contain more or fewer objects in total.
 1543 Figure A-2c visualizes just the object difference showing that children’s scenes contain
 1544 roughly as many objects on average as the initial example while adults’ scenes contain
 1545 around 0.75 fewer objects than are present in the initial example (dark shading in top row).
 1546 Thus, we opted to combine b and c by weighting the unsigned cone difference by the mean
 1547 inter-object distance across all comparisons to give our combined distance measure
 1548 (Figure A-2d and Figure 6c in the main text).

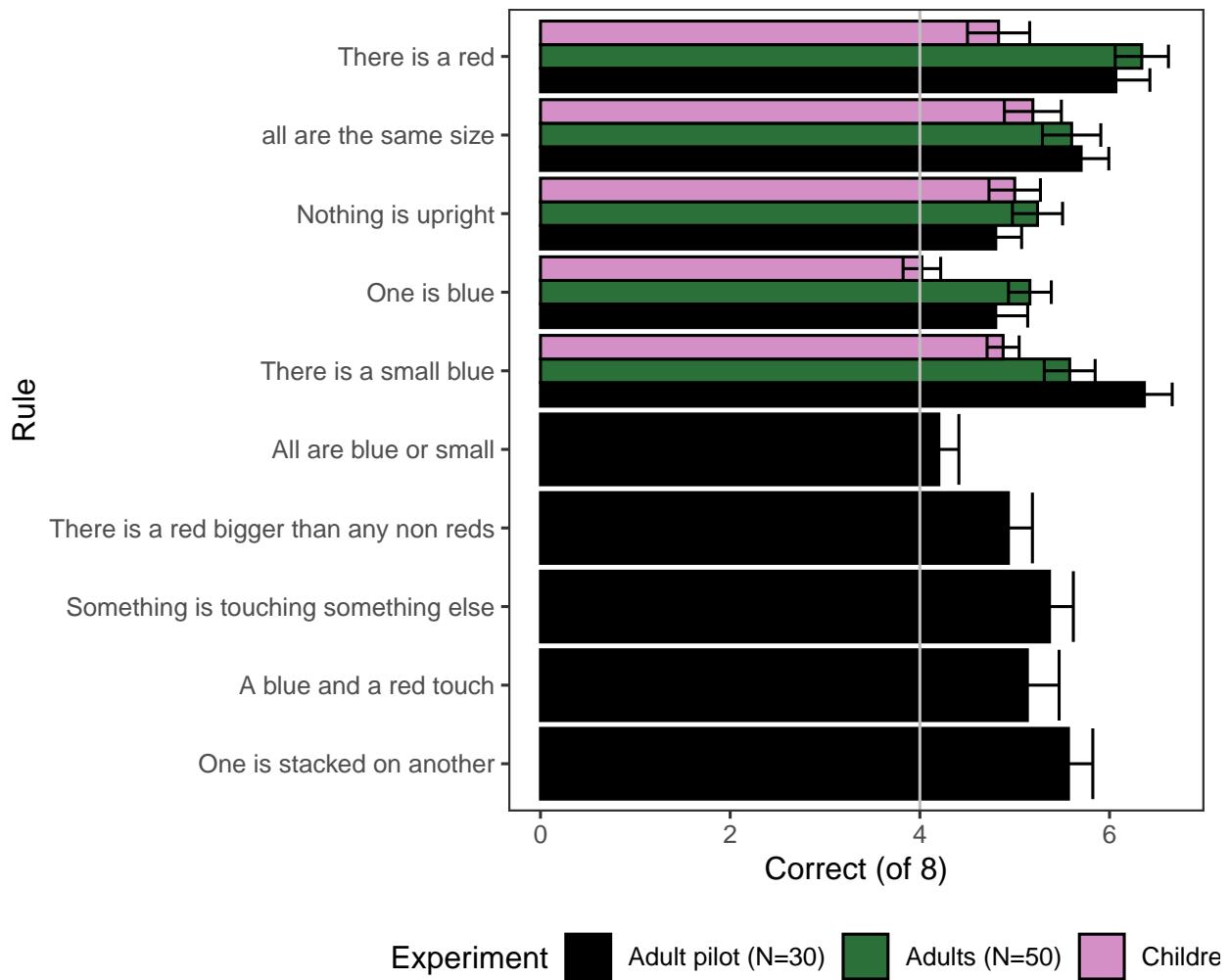
**Figure A-2**

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.

1549

Appendix D: Comparison with Bramley et al (2018)

1550 Finally, for interest and to demonstrate replication of our core results. We provide a
 1551 direct comparison between the generalization accuracies in the current sample of children
 1552 and adults and those in the sample of 30 adults modelled in (Bramley et al., 2018).
 1553 Bramley et al (2018) included 10 ground truth concepts, and the current paper uses just
 1554 the first five of these. Figure A-3 shows these accuracy patterns side by side, revealing the
 1555 adults in the current experiment performed approximately as well as those in the original
 1556 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).