

1 Running head: Teaching abstractions

2 **Teaching and learning generalizable abstractions**

3 August 28, 2024

4 Huang Ham¹, Bonan Zhao², Thomas L. Griffiths^{1,2}, Natalia Vélez¹

5 ¹Department of Psychology, Princeton University

6 ²Department of Computer Science, Princeton University

7 Keywords: pedagogy; social learning; Bayesian computational models; cultural evolution

8 Address for correspondence:

9 Natalia Vélez

10 Department of Psychology

11 Princeton University

12 Princeton, NJ 08540

13 E-mail: nvelez@princeton.edu

14

Abstract

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

A hallmark of effective teaching is that it grants learners not just a collection of facts about the world, but also a toolkit of abstractions that can be applied to solve new problems. How do humans transmit and acquire generalizable abstractions from examples? Here, we applied Bayesian models of pedagogy to a necklace-building task where teachers create necklaces to teach a learner “motifs” that can be flexibly recombined to create new necklaces. In Experiment 1 ($N = 151$), we find that human teachers produce necklaces that are simpler (i.e., have lower algorithmic complexity) than would be expected by chance, as indexed by a model that samples uniformly from all necklaces that contain the target motifs. This tendency to select simpler examples is partially captured by a pedagogical sampling model that tries to maximize the learner’s belief in the underlying motifs. In Experiment 2 ($N = 295$), we find that simplicity is beneficial. Human learners recover the underlying motifs better when teachers produce simpler sequences, and they learn best from human teachers rather than from model-generated examples. Our results suggest that the computational principles that underlie effective communication and teaching may also provide a first step towards understanding the transmission of culturally-specific abstractions.

1 Introduction

Teaching is a powerful means of transmitting culturally-specific abstractions, thereby laying the foundations to accumulate generalizable skills and knowledge across generations (Kline, 2015; Legare, 2019). One important kind of abstraction is a *motif*, a recurring pattern that can be composed into a larger work. For example, the basic stitches in knitting (knits and purls) are used to make recurring motifs or stitch patterns (e.g., stockinette, rib stitch) that can be flexibly combined to make many products (e.g., hats, scarves, sweaters). Motifs are found in many cultural products and often bear traces of the communities that produced them (Pesowski, Quy, Lee, & Schachner, 2020; Schachner, Brady, Oro, & Lee, 2018), such as the elaborate cable-knit patterns of Aran sweaters, the fundamental rhythmic pattern or *clave* of salsa music, and meander designs on the borders of Greek pottery.

Motifs pose a particular challenge for existing computational theories of social learning. Existing pedagogical sampling models characterize teaching and social learning as a series of recursive inferences: Teachers select examples that will maximize a learner’s belief in a target concept, and learners work backwards from the examples provided to infer the concept that the teacher is trying to communicate to them (Gweon, 2021; Shafto, Goodman, & Griffiths, 2014; Shafto, Wang, & Wang, 2021). This basic principle can explain a wide variety of communicative behaviors, including how teaching through demonstration differs from goal-directed behavior (Ho, Littman, MacGlashan, Cushman, & Austerweil, 2016; Tominaga, Knoblich, & Sebanz, 2022), how parents tune their speech to teach phonetic structures to infants (Eaves, Feldman, Griffiths, & Shafto, 2016), and how teachers improve their teaching based on feedback from learners (Chen, Palacci, Vélez, Hawkins, & Gershman, 2024). In addition, these computations appear to be neurally instantiated in mentalizing regions when teachers make decisions about what information to communicate to a learner (Vélez, Chen, Burke, Cushman, & Gershman, 2023).

However, prior work has largely focused on capturing how learners acquire solutions to particular problems (such as how to operate a particular toy; Aboody, Velez-Ginorio, Santos, &

58 Jara-Ettinger, 2023; Bridgers, Jara-Ettinger, & Gweon, 2020; Buchsbaum, Gopnik, Griffiths, &
59 Shafto, 2011) or identify the extension of particular categories (such as inferring the extent of a
60 hidden shape on a canvas; Shafto et al., 2014; Vélez et al., 2023). Most of these problems involve
61 a teacher choosing examples that constitute a part of the target that they intend to teach, such as a
62 single function of a toy with many functions or a single pixel inside a larger shape. Teaching
63 motifs presents the converse problem. For example, suppose an expert knitter draws a novice’s
64 attention to the rib stitch pattern on the collar of a sweater. Her goal is not to help the novice
65 identify sweaters based on this distinctive sub-component, as is the case in traditional teaching
66 games (Avrahami et al., 1997). Rather, the sub-component *itself* is the target of teaching; the
67 novice can later abstract away this sub-component to knit sock cuffs and hatbands. We do not yet
68 know to what extent pedagogical sampling models capture human behavior in this kind of
69 teaching problem, where examples of the whole are used to teach parts.

70 To approach this question, we studied how people teach and learn motifs within a simple
71 necklace-building task inspired by prior studies of cultural transmission (Clegg & Legare, 2016a,
72 2016b; Kleiman-Weiner et al., 2020). In Experiment 1 (N = 151), we tested to what extent
73 existing pedagogical sampling models capture how teachers transmit motifs. In this task, teachers
74 demonstrated motifs—recurring patterns of beads—by providing a single sample necklace.
75 Overall, our pedagogical sampling model provides a better quantitative fit to teachers’ decisions,
76 compared to a baseline model that samples uniformly from all necklaces that contain the target
77 motifs. The pedagogical sampling model also captures an important qualitative pattern in
78 teachers’ behavior: Teachers produce simpler examples than would be expected by chance. Next,
79 in Experiment 2 (N = 295), we tested the limits of existing models by giving human learners a
80 single sample necklace and asking them to both infer the underlying motifs and produce two new
81 necklaces that incorporate these motifs. Overall, learners perform better at both tasks when given
82 examples that are generated by human teachers or sampled from the pedagogical sampling model,
83 compared to necklaces generated by the baseline model. The effectiveness of these examples is
84 largely explained by simplicity—across the board, learners are better able to recover underlying

85 motifs when given simpler examples. However, learners performed best overall when given
86 human-generated examples, which suggests that state-of-the-art pedagogy models still miss
87 aspects of what makes human teaching effective. We close by discussing how pedagogical
88 sampling models could be extended to better capture how abstractions are culturally transmitted.
89 All experiment materials, data, and analysis code are publicly available at
90 https://osf.io/rnb9e/?view_only=099dad807964263a8e1196ce3dd2311.

91 **2 Computational framework**

92 **2.1 Task setup**

93 The experiments below use a necklace-building task as a simple case study of how people acquire
94 and transmit culturally-specific motifs. In prior work, necklace-building tasks have been used to
95 study basic mechanisms that drive cultural transmission, including how faithfully children imitate
96 (Clegg & Legare, 2016a, 2016b) and how adults learn concepts from step-by-step demonstrations
97 (Kleiman-Weiner et al., 2020). Experimental stimuli were adapted from Kleiman-Weiner et al.
98 (2020).

99 In Experiment 1 (Teaching Abstractions), participants play the role of expert artisans; their
100 task is to travel from one village to another, teaching an apprentice how to produce necklaces that
101 will sell well in each village. In our setting, each necklace is a string of 10 orange and green beads
102 that can be represented as a binary sequence. Each village has three favorite motifs, which are sub-
103 sequences of 2 or 3 beads that can be recombined to make necklaces (Figure 1A). A necklace sells
104 well in a village if and only if it includes all three of the motifs favored in that village (Figure 1B).
105 In Experiment 2 (Learning Abstractions), participants played the role of the apprentice. In each
106 village, participants saw one necklace generated by an expert teacher; their task was to infer the
107 underlying motifs in the necklace and to use these motifs to create new necklaces of their own.

108 In both experiments, we modeled teachers' and learners' behavior using two models: a
109 baseline model that assumes that the teacher samples uniformly among all necklaces that contain
110 the correct motifs, and a Bayesian pedagogy model that selects necklaces to teach that will

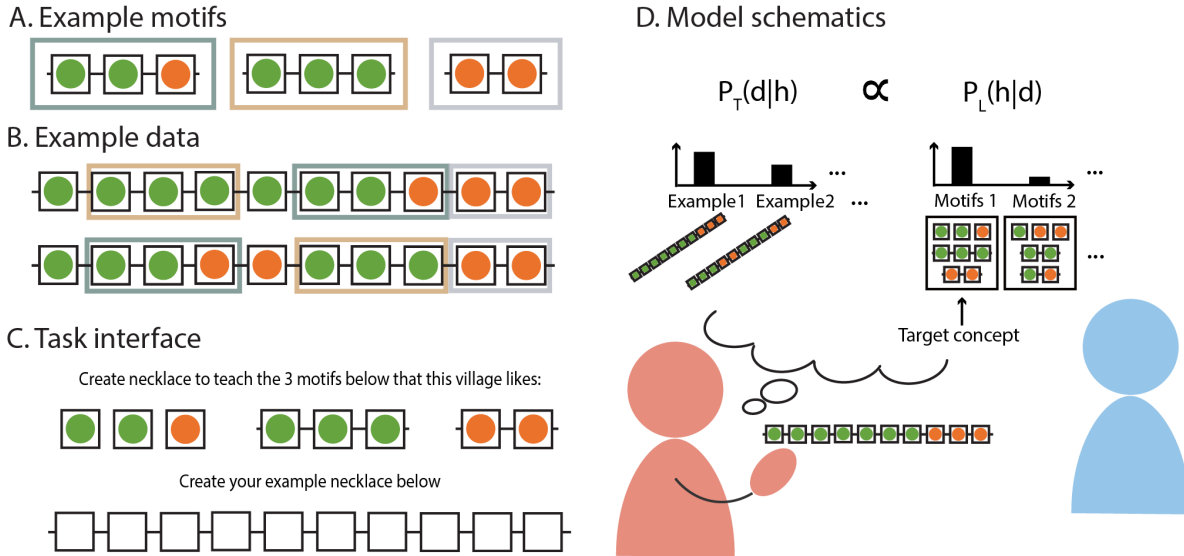


Figure 1: Teaching task. A-C. Experiment interface: A. On each trial, participants were shown the motifs favored by each village. B. These motifs can be recombined to create many new necklaces. C. Participants created a single sample necklace to teach learners these three motifs. D. Model schematic: The pedagogical sampling model actively selects necklaces d to teach ($P_T(d|h)$) by anticipating how learners will recover the underlying motifs h from the sample necklace ($P_L(h|d)$).

111 maximize the learner’s posterior beliefs in the correct motifs. To borrow terminology from prior
 112 work (Shafto et al., 2014), we will refer to these models as the “strong sampling” and
 113 “pedagogical sampling” models, respectively. We explain these models in more detail below.

114 2.2 Strong sampling model

115 Strong sampling refers to a process where examples are uniformly sampled from a target
 116 hypothesis (Shafto et al., 2014; Tenenbaum & Griffiths, 2001). In our task, each “example” is a
 117 sample necklace, and each “hypothesis” is a set of three motifs. In other words, the strong
 118 sampling model chooses uniformly among all necklaces that contain all three motifs favored by a
 119 particular village. We compared the strong sampling model to the behavior of both teachers and
 120 learners. In Experiment 1, comparing the fit of the pedagogical sampling model to that of the
 121 strong sampling model provides evidence about the extent to which teachers’ decisions are guided
 122 by higher-order inferences about a hypothetical learner’s beliefs (see *Teacher model*, below). In
 123 Experiment 2, we constructed a learner model that assumes that the teacher’s examples were

124 selected using strong sampling; we used this model as a baseline to measure how accurately
 125 human learners could be expected to recover the target hypothesis based on a single example (see
 126 *Learner baseline model*, below).

127 **2.2.1 Teacher model**

128 Given a hypothesis h , the strong sampling model predicts that the probability of selecting any
 129 sample necklace d is inversely proportional to the number of necklaces that contain the target
 130 motifs:

$$P_{\text{strong}}(d|h) = \begin{cases} \frac{1}{|h|} & \text{if } d \in h \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

131 We generated the predictions of the strong sampling model by first creating a matrix \mathbf{C} of
 132 hypotheses and data. The rows contain all possible 10-bead necklaces ($2^{10} = 1024$ possible
 133 necklaces) and the columns contain all possible triplets of motifs (with 2^2 possible 2-bead motifs
 134 and 2^3 3-bead motifs, there are $\binom{2^2+2^3}{3} = 220$ possible hypotheses). Each cell $c_{d,h}$ of this matrix
 135 indicates whether the necklace d could have been generated by combining the three motifs in h
 136 ($c_{d,h} = 1$) or not ($c_{d,h} = 0$). Thus, Equation 1 is equivalent to the following matrix operation:

$$P_{\text{strong}}(d|h) = \frac{c_{d,h}}{\sum_{d'} c_{d',h}}. \quad (2)$$

137 which corresponds to normalizing the entries of each column by its sum.

138 **2.2.2 Baseline learner model**

139 The baseline learner model assumes that the teacher selects a sample necklace d using strong
 140 sampling. We can use Bayes' rule to obtain the posterior distribution over motifs given the example
 141 provided:

$$P_{\text{baseline}}(h|d) = \frac{P_{\text{strong}}(d|h)P(h)}{\sum_{h'} P_{\text{strong}}(d|h')P(h')}. \quad (3)$$

142 where $P_{\text{strong}}(d|h)$ is the probability of selecting the sample necklace d under strong sampling when
 143 the true hypothesis is h (Equation 1) and $P(h)$ is the prior probability assigned to the hypothesis h .

144 The model further assumes a uniform prior over all sets of motifs (h), so $P(h)$ is the same for all h .
 145 Thus, Equation 3 can be simplified as:

$$P_{baseline}(h|d) = \frac{P_{strong}(d|h)}{\sum_{h'} P_{strong}(d|h')}. \quad (4)$$

146 **2.3 Pedagogical sampling model**

147 The pedagogical sampling model chooses necklaces to teach by anticipating how the learner will
 148 recover motifs from the example provided. Thus, rather than sampling uniformly from all
 149 necklaces that contain the target motifs, this model favors necklaces that are consistent with fewer
 150 alternative hypotheses. Intuitively, this strategy reduces the risk that learners will recover the
 151 wrong set of motifs from the sample necklace.

152 As an illustration, suppose that a village favors necklaces that contain the motifs “000”,
 153 “11”, and “001”. These motifs can be used to make many necklaces, including the following two
 154 examples:

$$\begin{aligned} 0000000111 &\rightarrow [000]00[001][11] \\ 0001100011 &\rightarrow 0[001]1[000][11] \end{aligned} \quad (5)$$

155 Here, the left side of each line shows the necklace as it would appear to a learner, and the right
 156 side breaks down each example into its component motifs. Note that both of these examples
 157 are ambiguous; there are several alternative sets of motifs that could have generated either of
 158 these necklaces. However, the model favors the necklace “0000000111” because there are fewer
 159 incorrect ways to parse it. Besides the 3 target motifs, this necklace is consistent with 4 incorrect
 160 motifs (i.e., “111”, “011”, “01”, “00”), which make up 13 consistent but incorrect alternative
 161 hypotheses (e.g., “00|000|11”, “11|01|000”). By contrast, “0001100011” can be parsed incorrectly
 162 in more ways. In addition to the 3 correct motifs, this necklace is consistent with 6 incorrect motifs
 163 (i.e., “00”, “01”, “10”, “100”, “011”, “110”), which make up 46 consistent but incorrect alternative
 164 hypotheses (e.g., “00|100|11”, “110|011|000”).

165 More formally, this model characterizes pedagogy as a form of cooperative communication
 166 between a teacher and a learner (Shafto et al., 2021). The model assumes that both the teacher
 167 and the learner have common knowledge about a space of hypotheses (i.e., all possible triplets of
 168 motifs) and a space of data (i.e., all possible necklaces; Shafto et al., 2014). The teacher selects
 169 necklaces to show to the learner that will maximize the learner’s beliefs in the true set of motifs
 170 favored by a particular village, and the learner works backwards from the necklace provided to
 171 infer what motifs the teacher is trying to demonstrate. These recursive inferences between the
 172 teacher and the learner are captured using the following system of equations:

$$P_{\text{learner}}(h|d) = \frac{P_{\text{teacher}}(d|h)p(h)}{\sum_{h'} P_{\text{teacher}}(d|h')p(h')}, \quad (6)$$

$$P_{\text{teacher}}(d|h) \propto P_{\text{learner}}(h|d)^\alpha, \quad (7)$$

174 where α is a free parameter that controls how strongly the teacher favors examples that maximize
 175 the learner’s posterior belief in the true motifs. In the results below, we fit α to individual
 176 participants’ responses in Experiment 1.

177 Following the procedure in Shafto et al. (2014), we calculated a solution to this system of
 178 equations using fixed-point iteration. We began by normalizing each column of the matrix \mathbf{C} by
 179 its sum, as in the strong sampling model (Equation 2). Next, we implemented the recursive
 180 inferences that distinguish pedagogical sampling from strong sampling by renormalizing this
 181 matrix. In the first iteration of this procedure, we normalized each row by its sum to generate
 182 $P_{\text{learner}}(h|d)$. Intuitively, each row of this matrix represents the posterior beliefs of a learner who
 183 assumes that the sample necklace d was generated using strong sampling. (Note that these first
 184 two iterations are equivalent to the strong sampling and learner models described in the previous
 185 section.) In the next iteration, we raised \mathbf{C} to the power of α and normalized each column by its
 186 sum to generate $P_{\text{teacher}}(d|h)$. Each column of this matrix represents the choice probabilities of a
 187 teacher who selects examples by considering the beliefs of the learner in the prior iteration.

188 Each iteration of this procedure represents an additional recursive inference. For example,

189 the next iteration yields the beliefs of a learner who tries to interpret what hypothesis the teacher
190 is attempting to communicate, and the next iteration after that yields the choice probabilities of a
191 teacher who tries to maximize the beliefs of a learner who actively interprets their examples, and
192 so on. In principle, we could iterate this system of equations indefinitely; in practice, $P_{\text{learner}}(h|d)$
193 and $P_{\text{teacher}}(d|h)$ converge to a fixed point after a finite number of steps (Wang, Wang, Paranamana,
194 & Shafto, 2020). We set the tolerance of convergence to 10^{-12} .

195 **3 Experiment 1: Teaching abstractions**

196 In Experiment 1, participants played the role of teachers. Their task was to provide a single
197 sample necklace to teach a learner generalizable “motifs” that could be recombined to produce
198 new necklaces that would sell well in a particular village. We compared the necklaces generated
199 by human teachers to those selected by a model that randomly selects a necklace that is consistent
200 with the target motifs (strong sampling) and a model that maximizes the learner’s belief in the
201 target motifs (pedagogical sampling).

202 **3.1 Methods**

203 **3.1.1 Participants**

204 We aimed to recruit 150 participants for the teaching task on Prolific (preregistration available
205 at https://aspredicted.org/GPT_HNV). In both experiments, participants were recruited through
206 the standard sample option; only participants who resided in the US, had an approval rate over
207 95%, were fluent in English, and completed 100 to 10000 studies were eligible to sign up. We
208 obtained data from 151 participants (M(SD) age = 39.78(13.73), 94 female, 53 male, and 4 non-
209 binary), potentially due to a server error at the time of submission. Participants earned \$3 for their
210 participation, and they were told that they could earn a performance bonus of up to \$1 based on
211 how well participants in Experiment 2 learned from the examples they selected. Thus, participants
212 were incentivized to provide helpful examples. In all experiments, participants provided informed
213 consent in accordance with the requirements of the Institutional Review Board.

214 3.1.2 Procedure

215 Participants played the role of master artisans; their task was to travel to 18 different villages and
216 teach an apprentice how to produce necklaces that would sell well in each village. On each trial,
217 participants were shown the motifs favored by a new village and were asked to create a single
218 10-bead necklace that contains all motifs favored by that village. Each village had a unique set of
219 motifs, and villages were presented in a randomized order.

220 On each trial, participants saw the 3 motifs favored by the village displayed on the top of
221 the screen, and they typed in a sample necklace by placing beads on an empty string in sequence.
222 For example, if a village had the motifs 000, 10, 111 teachers could pass on these motifs by
223 producing the necklace 0001011110 (emphasis added for demonstration). Participants could erase
224 beads from the sequence to correct mistakes and press a “submit” button when they were finished
225 with the sample necklace. To better align the task with our modeling assumptions, we constrained
226 participants’ responses so that they had to include all three motifs in their sample necklace; if the
227 necklace they produced was not valid, participants were prompted to correct the necklace before
228 proceeding.

229 3.1.3 Computational modeling

230 **Model fitting and comparison:** We compared models to participants’ responses using random-
231 effects Bayesian model selection (Rigoux, Stephan, Friston, & Daunizeau, 2014). First, we used
232 maximum likelihood estimation to fit the α parameter of the pedagogical sampling model to each
233 participant’s responses. For each participant, we then evaluated the fit of the strong sampling
234 and the pedagogical sampling models using the Bayesian information criterion ($BIC = k \log(n) -$
235 $2 \log(L)$), where k is the number of free parameters in the model (0 for strong sampling and 1
236 for pedagogical sampling), n is the number of trials observed per participant (which was fixed at
237 $n = 18$), and L is the maximized value of the probability of the data under the model. Finally, we
238 used $-0.5 \times BIC$ as an estimate of log model evidence for each participant (Bishop & Nasrabadi,
239 2006) and used it to compute the protected exceedance probability (pxp) for each model. Protected

240 exceedance probabilities treat models as random effects that can vary between participants; this
241 measure can be interpreted as the probability that a given model occurs most frequently in the
242 population.

243 **Model simulation:** In addition to the model comparison procedure described above, we also
244 used fitted models to create simulated datasets. Intuitively, simulated datasets allow us to compare
245 to what extent the necklaces produced by participants overlap with those that would be produced by
246 each model. To create simulated datasets, we first used the fitted α parameters for each participant
247 to create matrices of choice probabilities ($p_T(d|h)$) and then sampled from this matrix of choice
248 probabilities to obtain simulated responses.

249 **Measuring sequence complexity:** We measured the algorithmic complexity of sample
250 necklaces produced by human participants and in simulated datasets using the block
251 decomposition method (Zenil, Toscano, & Gauvrit, 2022). We chose this measure of algorithmic
252 complexity because of its theoretic connection with Kolmogorov complexity and algorithmic
253 information theory (Chaitin, 1969; Kolmogorov, 1965; Solomonoff, 1964) and also because it
254 captures well human subjective judgments of sequence randomness (Gauvrit, Singmann,
255 Soler-Toscano, & Zenil, 2016; Gauvrit, Zenil, Delahaye, & Soler-Toscano, 2014; Planton et al.,
256 2021). Intuitively, complexity scores provide an estimate of the length of the shortest program
257 needed to recreate each necklace.

258 **3.2 Results**

259 Both human participants and the pedagogical sampling model tended to create sample necklaces
260 that were simpler than those created by a strong sampling model (Figure 2A–B). We performed a
261 linear mixed-effects regression that predicted the algorithmic complexity of sample necklaces
262 based on fixed effects and random slopes of teacher type (human, pedagogical model, strong
263 sampling model) and random intercepts by villages (18 different ground-truth sets of motifs).
264 Both participants and the pedagogical sampling model created sample necklaces with lower
265 algorithmic complexity than those produced by the strong sampling model (human-generated vs.
266 strong-sampled necklaces: $b = -0.565, t(17.001) = -6.709, p < 0.001$; pedagogically-sampled

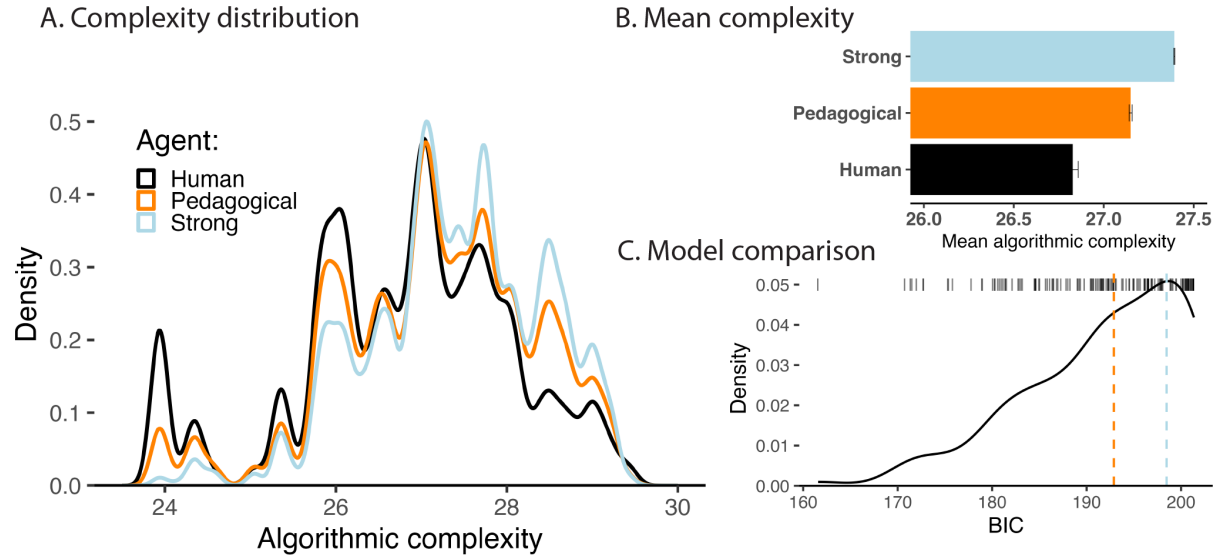


Figure 2: Experiment 1 results. A. Distribution of algorithmic complexity scores for human-generated necklaces (black line) and model-simulated necklaces (orange, blue lines). B. Average algorithmic complexity by teacher type. Human-created sample necklaces have the lowest mean complexity. Error bars denote standard error of the mean. C. Model comparison: Each tick-bar represents the fit between the pedagogical sampling model and a single participant’s responses, as measured by the Bayesian information criterion (BIC). Orange dotted line indicates the mean model evidences for the pedagogical sampling model. The blue dotted line indicates the BIC of the strong-sampling model; because this model selects uniformly among all valid necklaces, the probability of each participant’s responses under this model is a fixed value. Thus, points to the left of this blue line (lower BIC) are better fit by the pedagogical sampling model. Overall, participants’ responses were better captured by the pedagogical sampling model, compared to the strong sampling model.

267 vs. strong-sampled necklaces: $b = -0.243, t(17.003) = -7.108, p < 0.001$). Accordingly,
 268 participants’ choices were also better captured by the pedagogical sampling model ($p_{xp} = 1.000$;
 269 see Figure 2C for model evidences). These results suggest that the pedagogical sampling model
 270 captures a specific pattern in how people teach generalizable abstractions: Namely, effective
 271 teaching favors simpler examples. Note that we did not explicitly instruct the model to favor
 272 simpler examples—instead, this preference stems from a more general communicative principle.

273 However, the pedagogical sampling model alone does not account for the full pattern of
 274 participants’ responses. The examples generated by participants were still simpler than those
 275 simulated by the pedagogical sampling model (pedagogically-sampled vs. human examples:

276 $b = 0.322, t(16.997) = 5.648, p < 0.001$). Adding a penalty term for sequence complexity to the
277 pedagogical sampling model did not close this gap (see Figure S1). These results suggest that
278 participants' decisions may have been guided by additional inductive biases about what *kinds* of
279 simple examples are useful; if this is the case, then learners might derive benefits from
280 human-generated examples that are not fully captured by pedagogical sampling alone. In the
281 following experiment, we tested whether learners can indeed recover motifs from a single sample
282 necklace, and whether they learn best from examples generated by human teachers.

283 **4 Experiment 2: Learning abstractions**

284 Our results thus far suggest that reasoning about learners' mental states drives teachers to create
285 simpler necklaces than would be expected by chance. However, it is an open question whether
286 this simplicity aids learning—that is, whether the examples selected by teachers actually help
287 others recover the underlying motifs. In our second experiment, we approached these questions by
288 directly testing how well people learn triplets of motifs from observing a single sample necklace.
289 Sample necklaces were selected from those generated by human teachers in Experiment 1 and
290 from simulated datasets generated using the strong sampling and pedagogical sampling models.

291 **4.1 Methods**

292 **4.1.1 Participants**

293 We aimed to recruit 300 participants for the learning task on the Prolific platform using the standard
294 sample option (https://aspredicted.org/CK3_1NP). We lost data from 5 participants due to
295 server errors, leaving us with 295 participants for the learner task (M(SD) age = 40.02(12.31), 142
296 female, 149 male, and 4 non-binary). Participants were paid \$5 for completing the task, plus a
297 bonus of up to \$1 contingent on performance.

298 **4.1.2 Procedure**

299 Participants were told that they were apprentices to a master artisan. Their task was to travel to 18
300 villages and learn how to produce necklaces that would sell well in each village. As in Experiment

301 1, participants were told that necklaces would only sell well in a particular village if and only if
302 they contained all three motifs favored by that village. However, participants in Experiment 2 were
303 not shown these motifs directly; instead, they saw a single sample necklace generated by a teacher
304 and had to infer the motifs represented by the necklace.

305 On each trial, participants saw the sample necklace on the top of the screen and answered
306 two questions about the village’s motifs (Figure 3A). First, participants typed the motifs that they
307 believed that the village favors. As described above (*Computational framework*), each motif was
308 a sequence of 2 or 3 beads. Participants could erase the beads that they typed to correct mistakes,
309 and submit the motifs when they were ready. Once they submitted the three motifs, participants
310 could not change their responses. Next, participants were shown two empty necklaces and asked
311 to create two new necklaces that would sell well in that village. Participants could only proceed
312 if they produced two unique, length-10 necklaces that were distinct from the sample necklace.
313 Villages were presented in a random order.

314 Participants completed three within-subjects conditions, which differed in how sample
315 necklaces were generated. In the *Human* condition, participants saw sample necklaces created by
316 participants in Experiment 1. By contrast, in the *Pedagogical* and *Strong* conditions, participants
317 were shown sample necklaces that were simulated using the pedagogical- and strong-sampling
318 models, respectively. (See Experiment 1 for more details on how model-simulated necklaces were
319 generated.) Participants completed 18 trials total, comprising 6 trials for each teacher type.
320 Participants were blind to condition; that is, they did not know what type of teacher generated
321 each necklace. This experimental manipulation allows us to compare the effectiveness of human-
322 and model-generated examples.

323 **4.2 Results**

324 We measured participants’ performance using two outcome measures: The number of unique
325 motifs correctly reported by the participant (*correct motifs*; range: 0–3, where higher scores
326 indicate better performance) and the minimum number of changes that would have to be made to
327 change participants’ necklaces into a necklace that is consistent with the ground-truth motifs

A. Task interface

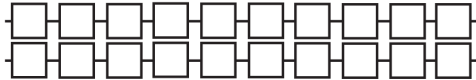
This is the example necklace created by your teacher to teach you the 3 motifs:



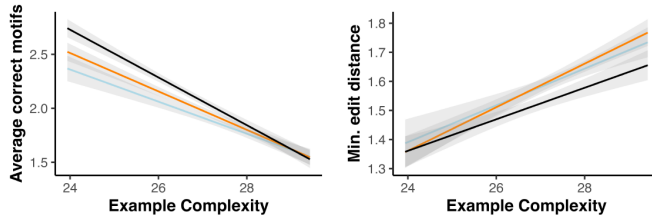
These are the 3 motifs you think the teacher tries to teach you:



Use the 3 motifs above to create 2 new necklaces to sell:



C. Effect of complexity on learning outcomes



B. Learning outcomes by teacher type

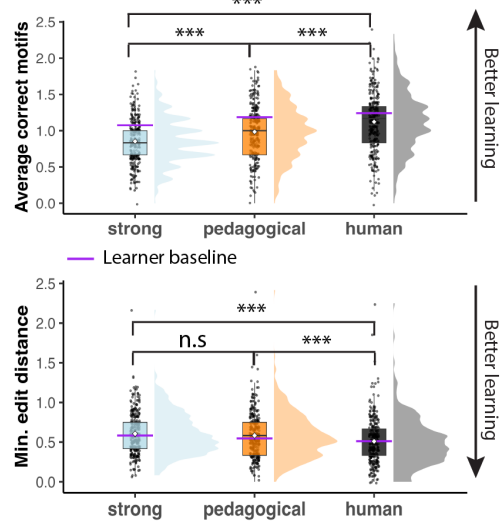


Figure 3: Experiment 2 results. A. Task interface: Participants saw a single sample necklace (top) that contained all three motifs favored by each village. They then explicitly reported the motifs (middle) and created two unique necklaces that were distinct from the sample necklace (bottom). B. Learning outcomes by teacher type: We measured learners’ performance based on the number of motifs that they correctly reported (range 1–3; higher scores indicate better performance) and the minimum number of edits needed to transform learners’ necklaces into a correct necklace (range 0–10; lower scores indicate better performance). Each point denotes a single participant’s average performance; white diamonds denote average scores by teacher type. As a baseline, we compared these scores to the performance of a learner model that infers the underlying motifs by assuming that teachers select examples using strong sampling (purple line). C. Correlation between the algorithmic complexity of the sample necklace and learner performance, as indexed by the number of correct motifs (left) and edit distance (right). Shaded areas denote 95% confidence intervals. Participants performed better when they were shown simpler sample necklaces.

328 (*minimum edit distance*; range: 0–10, where lower scores indicate better performance). We
 329 modeled each outcome using a mixed-effects ordinal regression with fixed effects of teacher type
 330 (i.e., Human, Pedagogical, Strong, with Strong as the reference level) and random slopes of
 331 teacher type by village.

332 Overall, participants learned best when they received sample necklaces selected by human
 333 teachers, rather than necklaces generated by the pedagogical- and strong-sampling models
 334 (Figure 3B). Participants reported more correct motifs when the examples were generated by the
 335 pedagogical-sampling model (Pedagogical vs. Strong: $b = 0.349, z = 5.120, p < 0.001$) and by
 336 human participants from Experiment 1 (Human vs. Strong: $b = 0.702, z = 8.011, p < 0.001$),

337 compared to the strong sampling model. However, participants recovered the most motifs overall
338 when they received examples generated by a human (Human vs. Pedagogical:
339 $b = 0.341, z = 4.230, p < 0.001$). Next, we found a similar pattern in the necklaces that learners
340 generated. Participants generated necklaces that were closer to the target motifs (i.e., had lower
341 minimum edit distances) when they received sample necklaces from a human teacher (Human vs.
342 Strong: $b = -0.336, z = -4.249, p < 0.001$), but not necklaces generated by the
343 pedagogical-sampling model (Pedagogical vs. Strong: $b = -0.091, z = -1.090, p = 0.276$).
344 Overall, participants produced more accurate necklaces when they received examples from a
345 human (Human vs. Pedagogical: $b = -0.244, z = -4.354, p < 0.001$).

346 On average, participants recovered approximately one of the three motifs specific to each
347 village (mean(SE) number of motifs: 0.987(0.015)) and produced necklaces that were less than
348 one bead away from an acceptable necklace (mean(SE) minimum edit distance: 0.566(0.013))
349 when provided a single sample necklace by a teacher. How well could participants be expected
350 to do, given the sparse and ambiguous information given to them? We benchmarked participants’
351 performance against the performance of the baseline learner model described above. We used
352 this model in two ways. First, we sampled a triplet of motifs from $P_{baseline}(h|d)$ (Equation 4) to
353 model how learners reported the motifs contained within each sample necklace. Second, the model
354 samples uniformly from all necklaces that contain this triplet to create two new necklaces.

355 Overall, we find that the learner baseline captures qualitative patterns in learner performance.
356 Like human learners, the baseline learner recovered approximately one of the three motifs specific
357 to each village (mean(SE) number of motifs: 1.168(0.004)) and produced necklaces that were less
358 than one bead away from an acceptable necklace (mean(SE) minimum edit distance: 0.547(0.006)).
359 To compare quantitative fits, we compared participants’ actual performance against this benchmark
360 using paired Wilcoxon tests. Regardless of teacher type, participants reported slightly *fewer* correct
361 motifs than the learner baseline (all $p < .05$ with the average difference of 0.178 motifs) and
362 produced necklaces with *similar* minimum edit distances as the learner baseline (all $p > .2$). (We
363 note one difference between the learner baseline and participants’ behavior: In 30% of trials,

364 participants reported motifs that were not consistent with the sample necklace they were provided.
365 Thus, it is possible that participants may have been inattentive. In an exploratory analysis, we
366 found that excluding these observations improved human learners' average performance slightly
367 but did not affect our interpretation of the results; see SI.) Thus, while participants underperformed
368 slightly relative to our baseline, they recovered about as much information as we could expect from
369 a single example.

370 To understand what makes human-generated examples particularly effective, we next used
371 mixed-effects ordinal regressions to model learners' performance based on an interaction between
372 teacher type and the algorithmic complexity of the sample necklaces provided; we also included
373 random effects of teacher type and algorithmic complexity by village. Participants who received
374 more complex sample necklaces reported fewer correct motifs (effect on correct motifs:
375 $b = -0.293, z = -3.699, p < 0.001$) and produced necklaces that were farther from correct
376 necklaces (effect on minimum edit distance: $b = 0.269, z = 4.721, p < 0.001$).

377 Together, our results suggest that simplicity is beneficial: Participants indeed learned better
378 when they were given simpler examples. However, participants still learned best when given
379 examples by human teachers—even though they were not aware of our experimental
380 manipulation—which suggests that existing models of teaching do not fully capture what makes
381 human teaching so effective. In the Discussion, we will consider how to bridge this gap.

382 **5 Discussion**

383 Teaching useful and generalizable abstractions underlies cultural and technological achievements
384 that require flexibility, innovation, and creativity. In this paper, we tested to what extent existing
385 models of teaching capture how humans teach generalizable abstractions. Overall, we found that
386 the general computational principles that underlie effective teaching and communication, as
387 formalized by the pedagogical sampling model, also capture a specific pattern in how humans
388 teach and acquire generalizable abstractions: Teachers favor simple examples, and learners learn
389 best from simple examples, without explicitly building simplicity as an assumption into our

390 models of either teachers or learners. However, our results also suggest that human teachers and
391 learners are even more sensitive to simplicity than this model would predict, highlighting an
392 exciting direction for future research.

393 Our results speak to prior theoretical debates on the importance of providing simpler data—
394 or “starting small”—for effective learning. “Starting small” refers to a hypothesis that learning
395 benefits from starting with simpler training data, both in humans and in artificial neural networks
396 (Elman, 1993; Rafferty & Griffiths, 2010; Zhao, Lucas, & Bramley, 2024). Our findings reveal
397 a similar phenomenon, where participants were better able to recover reusable abstractions when
398 they received simpler examples. Moreover, our model reveals that one reason why simplicity is
399 helpful is that it constrains the ways that learners can interpret the examples. However, the fact
400 that participants also favored examples that were *even simpler* than pedagogical sampling alone
401 would suggest that these models explain part but not all of what drives teachers to simplicity. We
402 speculate that there are additional inductive biases favoring simple patterns that recursive Bayesian
403 reasoning alone does not capture. For example, teachers may expect learners to parse necklaces
404 directionally, as though they were reading a script; for example, if learners “read” necklaces from
405 left to right, they may be more likely to interpret beads at the end of the necklace as overhang, rather
406 than as part of a motif. In addition, our model assumes that learners can perfectly evaluate whether
407 a necklace is consistent with a set of motifs. Human teachers may not share this assumption;
408 instead, they may favor even simpler necklaces to guide fallible learners who may make mistakes
409 when picking out motifs.

410 Overall, our findings provide a first demonstration of the usefulness of Bayesian pedagogy
411 for understanding how humans transmit generalizable, culturally-specific abstractions. However,
412 there are still many aspects of this domain that our simple task and model do not capture. Most
413 notably, our work examines how learners recover abstractions from a single data point. While our
414 findings show that learners can obtain some information even from this very sparse data, our task
415 stands in stark contrast to how skills are taught outside of the lab. Novices do not learn how to
416 knit a moss stitch or play musical chords from just a single example, but instead through repeated

417 interactions where an expert provides opportunities for learners to observe skills, corrects their
418 work, and sometimes provides explicit instruction (Kline, 2015). This additional structure has been
419 argued to be essential for the stable transmission of complex skills (Caldwell, Renner, & Atkinson,
420 2018; Tehrani & Riede, 2008). In addition, while the motifs in our task were simple sequences that
421 could be recombined in arbitrary ways, real-world cultural motifs are often imbued with meaning
422 (Cohn, 2012; Hawkins, Sano, Goodman, & Fan, 2023; Long, Fan, Huey, Chai, & Frank, 2024).
423 For example, skull motifs can remind viewers of the inevitability of death, and peonies appear
424 frequently in Chinese art as a symbol of prosperity and wealth. It is an open question how teachers
425 and learners coordinate on the meaning of motifs, or how these meanings constrain how motifs are
426 deployed. Thus, more work is needed to extend existing theories of pedagogy to fully embrace the
427 complexity of teaching generalizable abstractions.

428 Our work provides a theoretical and empirical framework to understand how teaching enables
429 learners to flexibly act, create, and innovate. We hope that revealing further components of this
430 ability will provide a fuller picture of how human intelligence is augmented by social learning and
431 culture.

References

- 432
- 433 Aboody, R., Velez-Ginorio, J., Santos, L. R., & Jara-Ettinger, J. (2023). When naïve pedagogy
434 breaks down: Adults rationally decide how to teach, but misrepresent learners' beliefs.
435 *Cognitive Science*, 47(3), e13257.
- 436 Avrahami, J., Kareev, Y., Bogot, Y., Caspi, R., Dunaevsky, S., & Lerner, S. (1997). Teaching by
437 examples: Implications for the process of category acquisition. *The Quarterly Journal of*
438 *Experimental Psychology Section A*, 50(3), 586–606.
- 439 Bishop, C. M., & Nasrabadi, N. M. (2006). *Pattern recognition and machine learning* (Vol. 4)
440 (No. 4). Springer.
- 441 Bridgers, S., Jara-Ettinger, J., & Gweon, H. (2020, Feb). Young children consider the expected
442 utility of others' learning to decide what to teach. *Nature Human Behaviour*, 4(2), 144–152.
443 doi: 10.1038/s41562-019-0748-6
- 444 Buchsbaum, D., Gopnik, A., Griffiths, T. L., & Shafto, P. (2011). Children's imitation of causal
445 action sequences is influenced by statistical and pedagogical evidence. *Cognition*, 120(3),
446 331–340.
- 447 Caldwell, C. A., Renner, E., & Atkinson, M. (2018). Human teaching and cumulative cultural
448 evolution. *Review of Philosophy and Psychology*, 9(4), 751–770. doi: 10.1007/s13164-017
449 -0346-3
- 450 Chaitin, G. J. (1969). On the length of programs for computing finite binary sequences: statistical
451 considerations. *Journal of the ACM (JACM)*, 16(1), 145–159.
- 452 Chen, A. M., Palacci, A., Vélez, N., Hawkins, R. D., & Gershman, S. J. (2024). A hierarchical
453 bayesian model of adaptive teaching. *Cognitive Science*, 48(7), e13477.
- 454 Clegg, J. M., & Legare, C. H. (2016a). A cross-cultural comparison of children's imitative
455 flexibility. *Developmental Psychology*, 52(9), 1435–1444.
- 456 Clegg, J. M., & Legare, C. H. (2016b). Instrumental and conventional interpretations of behavior
457 are associated with distinct outcomes in early childhood. *Child Development*, 87(2), 527–
458 542.
- 459 Cohn, N. (2012). Explaining “I can't draw”: Parallels between the structure and development of
460 language and drawing. *Human Development*, 55(4), 167–192.
- 461 Eaves, J., B. S., Feldman, N. H., Griffiths, T. L., & Shafto, P. (2016). Infant-directed speech
462 is consistent with teaching. *Psychological Review*, 123(6), 758–771. doi: 10.1037/
463 rev0000031
- 464 Elman, J. L. (1993). Learning and development in neural networks: The importance of starting
465 small. *Cognition*, 48(1), 71–99.
- 466 Gauvrit, N., Singmann, H., Soler-Toscano, F., & Zenil, H. (2016). Algorithmic complexity
467 for psychology: a user-friendly implementation of the coding theorem method. *Behavior*
468 *Research Methods*, 48, 314–329.
- 469 Gauvrit, N., Zenil, H., Delahaye, J.-P., & Soler-Toscano, F. (2014). Algorithmic complexity for
470 short binary strings applied to psychology: a primer. *Behavior Research Methods*, 46, 732–
471 744.
- 472 Gweon, H. (2021). Inferential social learning: Cognitive foundations of human social learning
473 and teaching. *Trends in Cognitive Sciences*, 25(10), 896–910.
- 474 Hawkins, R. D., Sano, M., Goodman, N. D., & Fan, J. E. (2023). Visual resemblance and
475 interaction history jointly constrain pictorial meaning. *Nature Communications*, 14(1),

476 2199.

477 Ho, M. K., Littman, M., MacGlashan, J., Cushman, F., & Austerweil, J. L. (2016). Showing
478 versus doing: Teaching by demonstration. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon,
479 & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 29). Curran
480 Associates, Inc.

481 Kleiman-Weiner, M., Sosa, F., Thompson, B., van Opheusden, B., Griffiths, T. L., Gershman, S.,
482 & Cushman, F. (2020). Downloading culture. zip: Social learning by program induction. In
483 *Proceedings of the 42nd annual meeting of the cognitive science society*.

484 Kline, M. A. (2015). How to learn about teaching: An evolutionary framework for the study of
485 teaching behavior in humans and other animals. *Behavioral and Brain Sciences*, 38, e31.
486 doi: 10.1017/S0140525X14000090

487 Kolmogorov, A. N. (1965). Three approaches to the quantitative definition of information.
488 *Problems of Information Transmission*, 1(1), 1–7.

489 Legare, C. H. (2019). The development of cumulative cultural learning. *Annual Review of*
490 *Developmental Psychology*, 1, 119–147.

491 Long, B., Fan, J. E., Huey, H., Chai, Z., & Frank, M. C. (2024). Parallel developmental
492 changes in children’s production and recognition of line drawings of visual concepts. *Nature*
493 *Communications*, 15(1), 1191.

494 Pesowski, M. L., Quy, A. D., Lee, M., & Schachner, A. (2020). Children use inverse planning
495 to detect social transmission in design of artifacts. In *Proceedings of the 42nd annual*
496 *conference of the cognitive science society*.

497 Planton, S., van Kerkoerle, T., Abbi, L., Maheu, M., Meyniel, F., Sigman, M., . . . Dehaene, S.
498 (2021, Jan). A theory of memory for binary sequences: Evidence for a mental compression
499 algorithm in humans. *PLoS Computational Biology*, 17(1), e1008598. doi: 10.1371/journal
500 .pcbi.1008598

501 Rafferty, A., & Griffiths, T. (2010). Optimal language learning: The importance of starting
502 representative. In *Proceedings of the 32nd annual meeting of the cognitive science society*.

503 Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection for
504 group studies—revisited. *Neuroimage*, 84, 971–985.

505 Schachner, A., Brady, T., Oro, K., & Lee, M. (2018). Intuitive archeology: Detecting social
506 transmission in the design of artifacts. In *Proceedings of the 40th annual meeting of the*
507 *cognitive science society*.

508 Shafto, P., Goodman, N. D., & Griffiths, T. L. (2014). A rational account of pedagogical reasoning:
509 Teaching by, and learning from, examples. *Cognitive Psychology*, 71, 55–89. doi: [https://](https://doi.org/10.1016/j.cogpsych.2013.12.004)
510 doi.org/10.1016/j.cogpsych.2013.12.004

511 Shafto, P., Wang, J., & Wang, P. (2021). Cooperative communication as belief transport. *Trends*
512 *in Cognitive Sciences*, 25(10), 826–828. doi: 10.1016/j.tics.2021.07.012

513 Solmonoff, R. (1964). A formal theory of inductive inference. Part I. *Information and Control*, 7,
514 224–254.

515 Tehrani, J. J., & Riede, F. (2008). Towards an archaeology of pedagogy: Learning, teaching and
516 the generation of material culture traditions. *World Archaeology*, 40(3), 316–331.

517 Tenenbaum, J. B., & Griffiths, T. L. (2001). Generalization, similarity, and bayesian inference.
518 *Behavioral and Brain Sciences*, 24(4), 629–640.

519 Tominaga, A., Knoblich, G., & Sebanz, N. (2022). Expert pianists make specific exaggerations
520 for teaching. *Scientific Reports*, 12, 21296. doi: 10.1038/s41598-022-25711-3

- 521 Vélez, N., Chen, A. M., Burke, T., Cushman, F. A., & Gershman, S. J. (2023). Teachers recruit
522 mentalizing regions to represent learners' beliefs. *Proceedings of the National Academy of*
523 *Sciences*, *120*(22), e2215015120. doi: 10.1073/pnas.2215015120
- 524 Wang, P., Wang, J., Paranamana, P., & Shafto, P. (2020). A mathematical theory of cooperative
525 communication. *Advances in Neural Information Processing Systems*, *33*, 17582–17593.
- 526 Zenil, H., Toscano, F. S., & Gauvrit, N. (2022). The block decomposition method. In *Methods*
527 *and applications of algorithmic complexity: Beyond statistical lossless compression* (pp.
528 125–149). Springer. doi: 10.1007/978-3-662-64985-5_6
- 529 Zhao, B., Lucas, C. G., & Bramley, N. R. (2024). A model of conceptual bootstrapping in human
530 cognition. *Nature Human Behaviour*, *8*(1), 125–136.

Supplemental Material

August 28, 2024

1 Building a complexity penalty into the pedagogical sampling model

In the main text, we report that teachers select examples that are even simpler than those selected by the pedagogical sampling model. As an exploratory analysis, we attempted to bridge this gap by extending the model with a penalty for more complex examples. We first obtained the learner’s posterior beliefs from the pedagogical sampling model, as described in the main text:

$$P_{\text{learner}}(h|d) = \frac{P_{\text{teacher}}(d|h)p(h)}{\sum_{h'} P_{\text{teacher}}(d|h')p(h')}, \quad (1)$$

Next, we defined the utility of teaching hypothesis h (a triplet of motifs) with example d (a sample necklace) by combining the learner’s posterior beliefs with the simplicity score:

$$U(d|h) = \ln(P_{\text{learner}}(h|d)) - wC(d) \quad (2)$$

where the negative surprisal term, $\ln(P_{\text{learner}}(h|d))$, captures the informational value of the example to the learner, and the complexity penalty, $C(d)$, is defined as the algorithmic complexity of d . w is the weight of the complexity penalty. Lastly, the utility score is converted into the probability of choosing d through a softmax function:

$$P_{\text{teacher}}(d|h) = \frac{e^{U(d|h)}}{\sum_{d'} e^{U(d'|h)}} \quad (3)$$

1.1 Results

We used the model estimation and model comparison procedures described in Experiment 1 of the main text. First, we used maximum likelihood estimation to fit the α and w parameters of the model. Next, we used Bayesian model selection to compare the fit of complexity penalty model (Figure S1, “with simplicity”) to the original pedagogical sampling model (Figure S1, “without simplicity”). Even after directly penalizing complex examples, the original pedagogical sampling model best captures the behavior of human teachers (PXP = 1). These results suggest that people may not favor simpler examples for simplicity’s sake;

instead, the decisions of human teachers may be guided by additional inductive biases that are not captured by existing theories. We return to this point in the General Discussion.

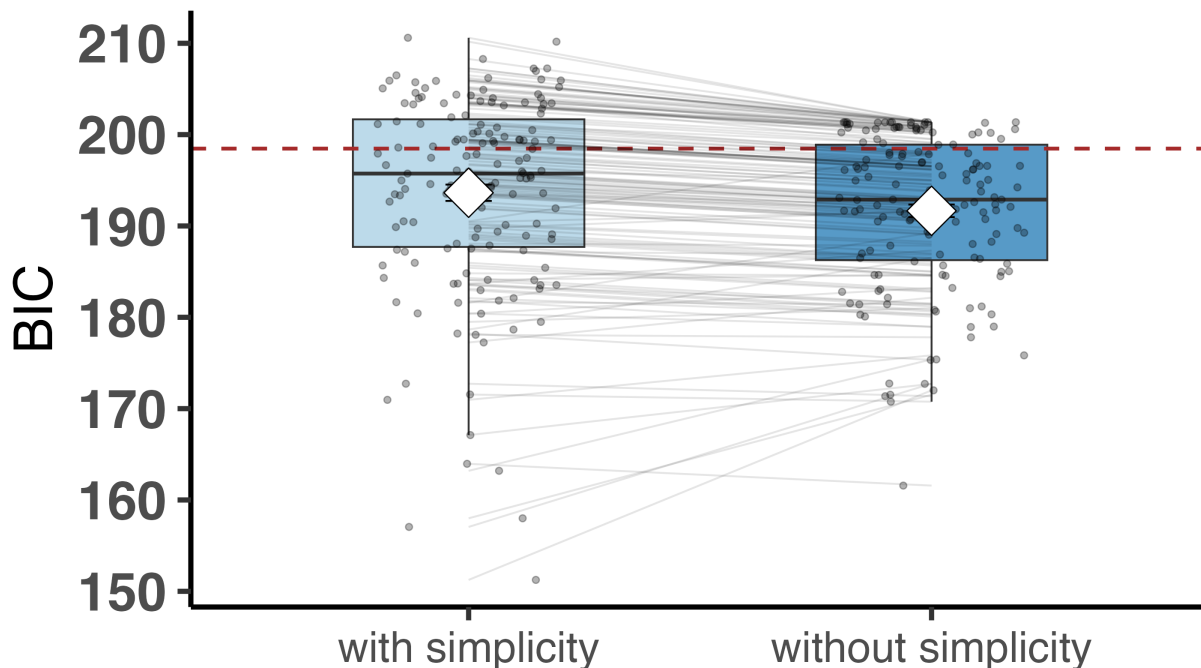


Figure S1: Bayesian information criterion (BIC): each dot represents the BIC of the model of each participant. The red dotted line indicates the BIC of the strong-sampling model. This shows overall, the original pedagogical model fit better (smaller BIC) than the pedagogical model that includes the simplicity score.

2 Stimuli

We chose 18 triplets of motifs for the 18 villages that participants visit in both the teacher task (Experiment 1) and learner task (Experiment 2). 9 villages favor triplets of motifs that contain 2 motifs with 2 beads and 1 motif with 3 beads, while the remaining villages favor triplets of motifs that contain 1 motif with 2 beads and 2 motifs with 3 beads. We avoided triplets of motifs that contain either all length-2 or length-3 motifs because the pedagogical sampling and strong sampling model do not make clearly distinguishable predictions for this subset of stimuli. The stimuli are: “10|010|011”, “01|011|101”, “01|10|010”, “00|11|010”, “00|01|110”, “01|11|101”, “10|11|110”, “00|011|100”, “11|010|101”, “01|101|110”, “00|10|011”, “10|010|101”, “11|001|011”, “01|11|010”, “01|10|101”, “00|100|101”, “00|01|010”, “01|010|101”. “1” represents orange beads and “0” represents green beads.

3 Comparing learner performance to the baseline learner model after exclusions

In the main text, we note that human learners sometimes reported motifs that were inconsistent with the sample necklace they had received. It is possible that these trials reflect instances where participants were inattentive or made typing errors. Therefore, as an exploratory analysis, we also compared human learners' performance to the baseline learner model after excluding these trials. After excluding the 30% of the trials where the inferred motifs were not consistent with the example necklace, we overall saw a slight improvement in learners' performance compared to the learner baseline model. On average, participants recovered approximately one of the three motifs specific to each village (mean(SE) number of motifs: 1.091(0.015)) and produced necklaces that were less than one bead away from an acceptable necklace (mean(SE) minimum edit distance: 0.506(0.011)). The baseline learner recovered approximately one of the three motifs specific to each village (mean(SE) number of motifs: 1.163(0.007)) and produced necklaces that were less than one bead away from an acceptable necklace (mean(SE) minimum edit distance: 0.550(0.009)). After exclusion, participants still reported *fewer* correct motifs than the learner baseline if the examples were chosen by the models (all $p < .05$). However, if examples came from human teachers, participants reported a *similar amount* of correct motifs as the baseline learner model ($p = .80$). Participants also produced necklaces with *similar* minimum edit distances as the learner baseline if the examples came from the pedagogical sampling model ($p = .285$). However, if the examples came from the strong sampling model or human teachers, participants produced necklaces with *smaller* minimum edit distances than the learner baseline (all $p < .01$).