

Learning generalizable representations through efficient coding

Zeming Fang^{1,2} and Chris Sims^{2*}

¹ Brain Health Institute, National Center for Mental Disorders, Shanghai Mental Health Center, Shanghai Jiao Tong University School of Medicine and School of Psychology, Shanghai, 200030, China

⁹ ²Department of Cognitive Science, Rensselaer Polytechnic Institute, Troy, NY, United States

10 * Corresponding author.

12 Author Note

13 Zeming Fang  <https://orcid.org/0000-0002-8091-4413>

14 Chris Sims <https://orcid.org/0000-0002-3110-1686>

16 **Abstract**

17 Reinforcement learning (RL) theory explains human behavior as driven by the goal of maximizing
18 reward. However, conventional RL approaches offer limited insights into human generalization.
19 Here, we propose refining the classical RL framework by incorporating an efficient coding
20 principle, which emphasizes maximizing reward using the simplest necessary representations.
21 This refined framework predicts that intelligent agents, constrained by simpler representations,
22 will inevitably: 1) distill environmental stimuli into fewer, abstract internal states, and 2) detect
23 and utilize rewarding environmental features. Consequently, complex stimuli are mapped to
24 compact representations, forming the foundation for generalization. We tested this model in two
25 experiments designed to reveal human generalization abilities. Our results demonstrate that
26 while classical RL models, which do not simplify representations, fail to generalize, an efficient
27 coding model that learns compact representations successfully achieves human-level
28 performance. We argue that the classical RL objective, augmented with efficient coding,
29 represents a more comprehensive computational framework for understanding human behavior
30 in both learning and generalization.

31

32

1. Introduction

33 The aphorism “A man can never step into the same river twice” speaks to the ever-changing
34 nature of the world. Making sense of this dynamic reality requires the ability to *generalize*; that
35 is, to extract knowledge from past experiences and apply it to new, unseen futures. Effective
36 generalization remarkably improves the capacity of intelligent agents to adapt to rapid changes.
37 For example, consider a child learning to ride a bike. She makes numerous attempts, falling and
38 adjusting her balance through trial and error. Once bike riding is mastered, the child can then
39 generalize those balancing skills to ride a scooter, allowing her to quickly master the scooter
40 without having to learn from scratch. Given its importance to adaptive learning, generalization
41 has been the focus of study in both cognitive neuroscience ([Shepard, 1987](#); [Shohamy & Wagner, 2008](#);
42 [Sims, 2018](#)) and machine learning ([Asadi et al., 2018](#); [Li et al., 2003](#); [Pensia et al., 2018](#)).

43 Recent research illustrates that representation learning is one of the cornerstones that
44 support generalization ([Bengio et al., 2013](#); [Goodfellow et al., 2016](#); [Radulescu et al., 2021](#)).
45 Representation learning involves the transformation of raw environmental stimuli or events into
46 robust abstract states (“*state abstraction*”), which summarize underlying patterns and
47 regularities in the raw data. For example, a bike and scooter may be conceptually abstracted into
48 transportation, enabling a child to realize they can transfer balancing skills previously learned
49 from riding a bicycle to a scooter. In addition, effective representations can detect and extract a
50 subset of the most informative and rewarding features within environments (“*rewarding feature*
51 *extraction*”). For instance, although bicycles and scooters have distinct designs, their shared
52 feature of having two wheels requires similar balancing skills. Historically, there has been a gap
53 in the theoretical and comprehensive understanding of how to constitute effective
54 representations. Bridging this gap and developing algorithms that learn generalizable
55 representations has become a central pursuit in recent research on human cognitive
56 neuroscience ([Flesch et al., 2022](#); [Nelli et al., 2023](#); [Niv, 2019](#); [Radulescu et al., 2021](#)) and artificial
57 intelligence ([Bengio et al., 2013](#); [Higgins et al., 2016](#); [Li et al., 2006](#); [Shwartz-Ziv, 2022](#); [Tishby et
58 al., 2000](#)).

59 This paper focuses on understanding how humans learn effective representations that
60 enhance their generalization abilities. One influential framework for understanding human

61 behavioral learning is reinforcement learning (RL), which views intelligent behavior as
62 subservient to the maximization of reward ([Silver et al., 2021](#); [Sutton & Barto, 2018](#)). This
63 framework provides a normative understanding of a spectrum of human learning processes ([Daw](#)
64 [et al., 2011](#); [Jiang et al., 2023](#); [Niv & Langdon, 2016](#); [Rescorla, 1972](#); [Ribas-Fernandes et al., 2011](#);
65 [Tomov et al., 2021](#); [van Opheusden et al., 2023](#); [Xia & Collins, 2021](#)) and offers theories on the
66 underlying neural mechanisms ([Barto, 1995](#); [Montague et al., 1996](#); [Niv, 2009](#); [Schultz et al.,](#)
67 [1997](#)). However, by itself, the traditional RL framework provides very limited insights into human
68 representation learning and generalization ([Gershman & Daw, 2017](#); [Ho et al., 2022](#); [Mnih et al.,](#)
69 [2015](#); [Niv, 2019](#); [Niv & Langdon, 2016](#)). The framework often assumes a predefined, fixed set of
70 task representations on which learning can operate directly, without the need for additional
71 representation learning ([Sutton & Barto, 2018](#)). However, in real-world decision-making, humans
72 are not provided with predefined representations. Instead, they must infer these representations
73 from complex and dynamic environmental observations.

74 Here, we propose augmenting the classical RL theory to incorporate the principle of
75 efficient coding ([Barlow, 1961](#)): while maximizing reward, intelligent agents should use the
76 simplest necessary representations. Our modification to the RL framework is not unfounded
77 because the human brain, as a biological information processing system, possesses finite
78 cognitive resources ([Miller, 1956](#)). The idea of efficient use of cognitive resources has had
79 profound impacts across many domains in psychology and neuroscience, including perception
80 ([Simoncelli & Olshausen, 2001](#); [Sims, 2016](#); [Wei & Stocker, 2015](#)), working memory ([Bates et al.,](#)
81 [2019](#); [Sims et al., 2012](#)), perceptual-based generalization ([Sims, 2018](#)), and motor control ([Lerch](#)
82 [& Sims, 2021](#)). Furthermore, our approach aligns with Botvinick's ([2015](#)) proposal that the
83 efficient coding principle can be instrumental in understanding the representation of problems
84 in learning and decision-making.

85 Critically, our proposed approach suggests that, driven by the principle of efficient coding,
86 an intelligent agent can autonomously learn appropriate simplified representations, which
87 enables both state abstraction and the extraction of rewarding features, naturally resulting in
88 generalization. To validate these predictions, we designed two experiments focusing on learning
89 and generalization. Participants first learned a set of stimulus-action associations and were then

90 tested on their ability to generalize to a new set of associations they had not encountered before.
91 The first experiment investigates the emergence of state abstraction, while the second explores
92 the extraction of rewarding features. Human participants displayed strong generalization abilities
93 in both experiments, correctly responding to new associations without additional training. We
94 developed a principled model based on efficient coding and demonstrated its capacity to achieve
95 human-level generalization performance in both experiments—performance that classical RL
96 models have not accomplished. These findings lead us to conclude that generalization is an
97 inherent outcome of efficient coding. Given humans' remarkable capacity for generalization, we
98 assert that the classical RL objective augmented with efficient coding and reward maximization
99 presents a more comprehensive computational objective for human learning.

100

101 **2. Results**

102 Humans exhibit two types of generalizations: *perceptual-based* and *functional-based*
103 generalizations. Perceptual-based generalization occurs when two stimuli share a similar
104 appearance ([Shepard, 1987](#); [Sims, 2016](#); [Sims et al., 2012](#)). Functional-based generalization, in
105 contrast, occurs between stimuli that have similar functions (e.g. linked to the same actions),
106 even when they do not look alike ([Collins & Frank, 2013](#); [Collins & Frank, 2016](#); [Meeter et al.,](#)
107 [2009](#); [Myers et al., 2003](#); [Shohamy & Wagner, 2008](#)). The latter type of generalization is more
108 complex because it necessitates the acquisition of unseen environmental statistics before it can
109 occur.

110 To investigate both types of generalization, we leveraged the *acquired equivalence*
111 paradigm ([Meeter et al., 2009](#); [Myers et al., 2003](#); [Shohamy & Wagner, 2008](#)). This experimental
112 framework first links two visually distinct stimuli with identical actions, then assesses the increase
113 in generalization between these stimuli based on their shared actions. This approach effectively
114 establishes the functional similarity between the two stimuli, enabling a controlled experimental
115 investigation into participants' ability of functional-based generalization.

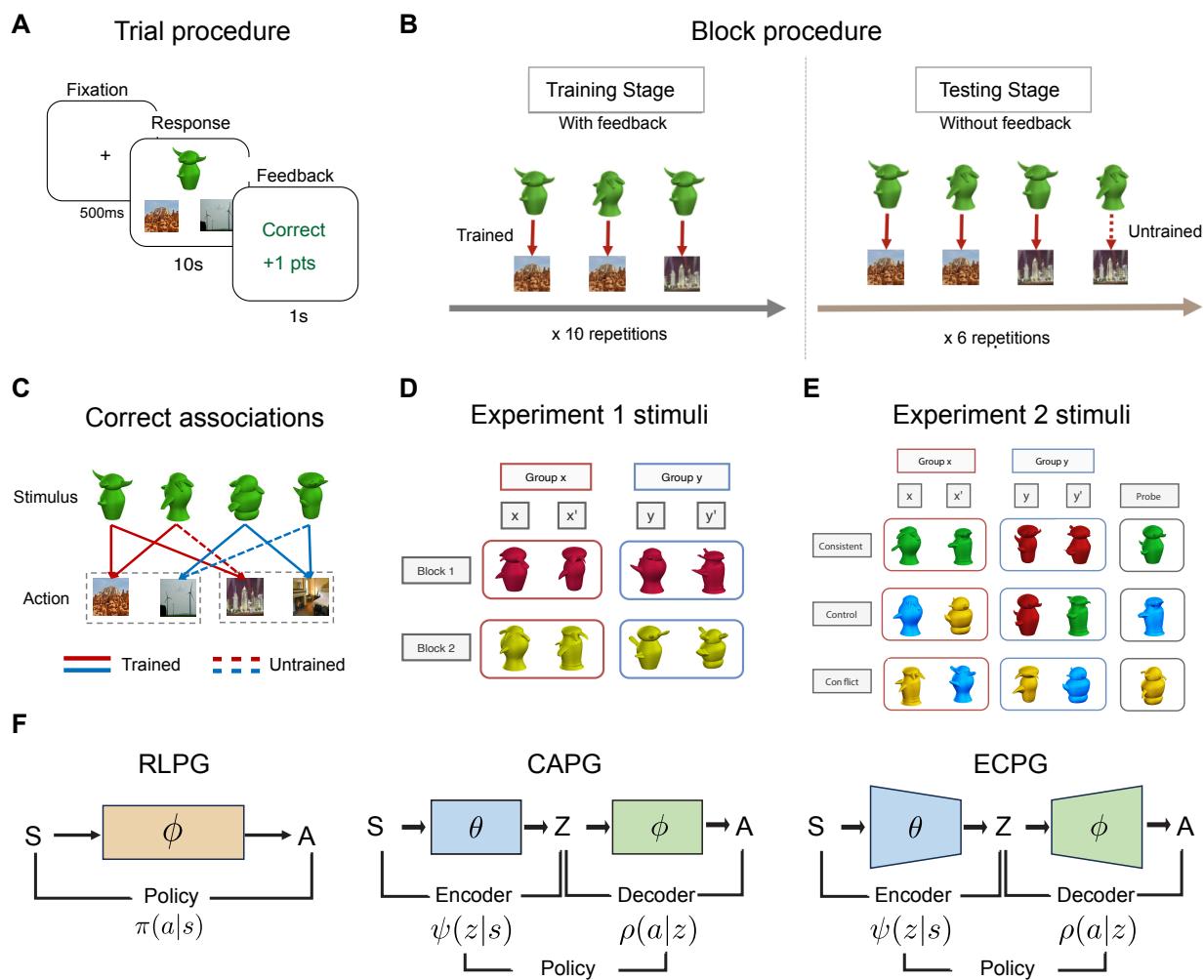
116 Specifically, participants performed a two-stage task. In each trial, participants were
117 shown an alien (stimulus s) and were told that different aliens preferred to visit different
118 locations. For a given stimulus, participants were required to choose one of two places (action a)
119 that they believed the alien would prefer to visit (**Fig. 1A**). During the training stage, participants
120 were trained on six stimulus-action associations, each repeated ten times to learn the
121 equivalence between stimuli based on their associated actions (**Figs. 1B and 1C**). For example, if
122 aliens s_1 and s_2 both preferred to visit a desert (a_1) rather than a forest (a_2), then they are
123 equivalent and the psychological similarity of the two aliens may increase. During the training
124 stage, participants received feedback (reward r , taking a value of either 0 or 1) after every choice.

125 In the testing stage, participants were tested on eight associations: the six trained
126 associations plus two untrained associations that were not presented in the training stage. The
127 untrained associations were used to evaluate people's generalization performance. For example,
128 if the participant learned during the training stage that s_1 and s_2 were similar to each other (had
129 similar preferences), then participants might generalize other preferences from s_1 to s_2 , even

130 though no feedback was given about those preferences. No feedback was provided during the
 131 testing stage, and each association was repeated six times.

132 To quantify human generalization ability, we calculate the “untrained accuracy”, which is
 133 the response accuracy for the untrained associations that were not presented during training.
 134 The higher the untrained accuracy, the better a participant’s generalization ability. Similarly,
 135 “trained accuracy”—the response accuracy for trained associations that were presented in the
 136 training stage—serves as a measure of human learning performance. Both metrics are crucial and
 137 will be used extensively throughout this paper.

138 All data were collected online via Amazon Mechanical Turk.



139

140 **Figure 1** Acquired equivalence experimental paradigm and model architectures.

- 141 **A.** One trial consists of three screens: a 500-ms fixation screen, a 10-s response screen, and
142 a 1-s feedback screen. Each response screen displays an alien stimulus, as well as two
143 location pictures representing different actions.
144 **B.** One block contains two stages. The training stage trains three associations with feedback.
145 The testing stage tests an untrained association (dashed line) in addition to the three
146 trained associations without feedback.
147 **C.** One block contains two groups, each with two stimuli. The incorrect actions of one group
148 correspond to the correct actions of the other. The structure of the stimuli is not disclosed
149 to participants.
150 **D.** Stimuli used in Experiment 1 are designed to be the same color but with different shapes
151 and appendages to control for perceptual similarity. The four stimuli are referred to as
152 x, x', y, y' , with x and x' associating with the same actions, as do y and y' .
153 **E.** Stimuli used in Experiment 2. Each block contains a different type of perceptual similarity.
154 **F.** The model architectures. The classical reinforcement learning policy gradient (RLPG)
155 model learns a policy that maps from stimuli s to a distribution of action a . Due to the
156 introduction of representation r , the policies of the cascade policy gradient (CAPG) and
157 the efficient coding policy gradient (ECPG) model are broken into an encoder ψ and
158 decoder ρ . The ECPG and CAPG model have the same architecture except that the ECPG
159 model optimizes to use simpler representations z .

161 2.1 Modeling human behavior at the computational level

162 David Marr ([1982](#)) famously argued that the human brain can be understood at three levels: *the*
163 *computational level*, which defines the goals to be achieved; *the algorithmic level*, which details
164 the specific algorithms the human brain used to reach these goals; and *the implementational*
165 *level*, which describes how these algorithms are physically realized.

166 In psychology and cognitive science, researchers often build models at the algorithmic
167 level. They typically postulate specific cognitive mechanisms within the human brain, describe
168 these mechanisms using computer programs, and demonstrate that incorporating them provides
169 a better explanation of human behavioral data (e.g. [Collins & Frank, 2013](#); [Collins & Frank, 2016](#);
170 [Niv et al., 2015](#)).

171 However, the question of whether the human brain reconstructs efficient representations
172 for task stimuli is situated at the computational level. Therefore, we need to construct models at
173 this same level. In concrete, we formalized our hypotheses—with or without efficient coding—
174 as distinct computational goals, each addressed using the simplest possible algorithm. Crucially,
175 unlike the algorithmic-level models, the computational-level models do not presume any specific
176 mechanisms. Instead, these mechanisms naturally emerge during the process of achieving the
177 defined computational goal. Thus, computational-level models not only explain human behaviors

178 but also shed light on the potential cognitive mechanisms underlying these behaviors, thereby
179 demonstrating superior explanatory power over algorithmic-level models.

180 We built three computational-level models. First, we established a classical RL baseline,
181 named *Reinforcement Learning Policy Gradient* (RLPG; see **Fig. 1F** and **Method 4.4.1**), which
182 assumes that humans do not learn simplified representations. The computational goal is
183 formulated as follows:

184
$$E_{\pi}[r(s_t, a_t) - b] \quad (1)$$

185 where $\pi(a|s)$ is a policy that maps a stimulus, s , to a distribution of actions, a . On each
186 trial, an agent had to choose between two possible actions, each with a 50% chance of being
187 correct. Prior to making a decision, the agent was expected to have a baseline reward expectation
188 of $b = 0.5$. This baseline was used to evaluate the “goodness” of the actual reward received. A
189 reward was considered positive if it exceeded the agent's expectation, otherwise negative. The
190 RLPG model interpreted human behavior as involving the search for the policy that yielded the
191 greatest reward (above the baseline) in the process of interacting with the environment.

192 Second, we developed an *efficient coding policy gradient* model (ECPG; **Fig. 1F** and
193 **Method 4.4.2**), which posits that humans learn simpler representations through efficient coding.
194 The challenge in modeling this principle lies in defining the complexity (or simplicity) of
195 representations. Recent studies on human perception have conceptualized perception as an
196 information transmission process, where an encoder transmits environmental sensory signals (s)
197 into internal representation (z) ([Bates et al., 2019](#); [Sims, 2018](#); [Sims et al., 2012](#)). These studies
198 measure the complexity of representations by the amount of information transmitted by the
199 encoder, quantified by the mutual information between stimuli and representations $I^{\psi}(S; Z)$.
200 Based on these work, the computational goal of efficient coding is formalized as maximizing
201 reward while minimizing the representation complexity,

202
$$E_{\psi,\rho}[r(s_t, a_t) - b] - \lambda I^{\psi}(S; Z) \quad (2)$$

203 The critical parameter $\lambda \geq 0$, referred to as the *simplicity parameter*, controls for the
204 tradeoff between the classical RL objective and representation simplicity. When $\lambda = 0$, the agent
205 does not compress stimuli representations for simplicity, and the efficient coding goal reduces to
206 the RL goal. Conversely, as $\lambda \rightarrow \infty$, the agent learns the simplest set of representations,

207 encoding all stimuli into a single, identical representation. Therefore, the optimal λ should be a
208 moderate value, balancing compressing without oversimplification. Due to the introduction of
209 latent representation z , the policy needs to be broken down into an encoder, ψ , and a decoder,
210 ρ , which are simultaneously optimized according to Eq. 2 (**Fig. 1F**).

211 To test whether humans learn simpler representations, the establishment of the RLPG
212 and ECPG models would typically be sufficient, because the contrasting hypotheses they
213 represent (RLPG stands for “No”, ECPG stands for “Yes”) together cover the entire hypothesis
214 space. One concern, however, is that the introduction of the representation in ECPG has changed
215 the model architecture, potentially introducing confounding factors. To control these
216 confounders, we implemented a third model, *cascade policy gradient* (CAPG; Method 4.4.3),
217 which also supports the non-efficient coding hypothesis. The CAPG is a special case of the ECPG
218 model which sets the simplicity parameters to 0 ($\lambda = 0$ in Eq. 2) (**Fig. 1F**),

219
$$E_{\psi, \rho}[r(s_t, a_t) - b] \quad (3)$$

220 This model serves as an intermediary between the RLPG and ECPG models, optimizing for the
221 classical RL objective while concurrently updating the representations.

222 To ensure that observed behavioral differences result only from optimizing different
223 computational goals, we carefully controlled for all other model components to avoid any
224 confounding factors. First, all three models address their computational goals using the same
225 *policy gradient* approach, where models explicitly learn and maintain a parameterized policy
226 ([Sutton & Barto, 2018](#); [Williams, 1992](#)). The method was selected over the more commonly used
227 value function approach in psychology and neuroscience because it introduces a minimum
228 number of parameters, therefore better distilling the computational essence of each
229 computational goal. Second, the three models were initialized to (nearly) the same state. Due to
230 the distinct appearances of stimuli in the experiments, we pretrained the encoders of the CAPG
231 and ECPG models to achieve a 99% initial discrimination accuracy among the four stimuli (see
232 **Method 4.5**). We chose a threshold of 99% instead of 100% for two reasons: first, to model the
233 perceptual noise present in the human visual system, and second, to prevent gradient vanishing,
234 which is an engineering concern. The RLPG model implicitly assumes perfect discrimination
235 between stimuli and, therefore, does not require the same pretraining as others. Lastly, we used

236 the same model fitting method for all three models, fitting the parameters to each participant
237 separately using maximum-a-posteriori (MAP) estimation (**Method 4.8**), based on behavioral
238 data from both the training and the testing stages.

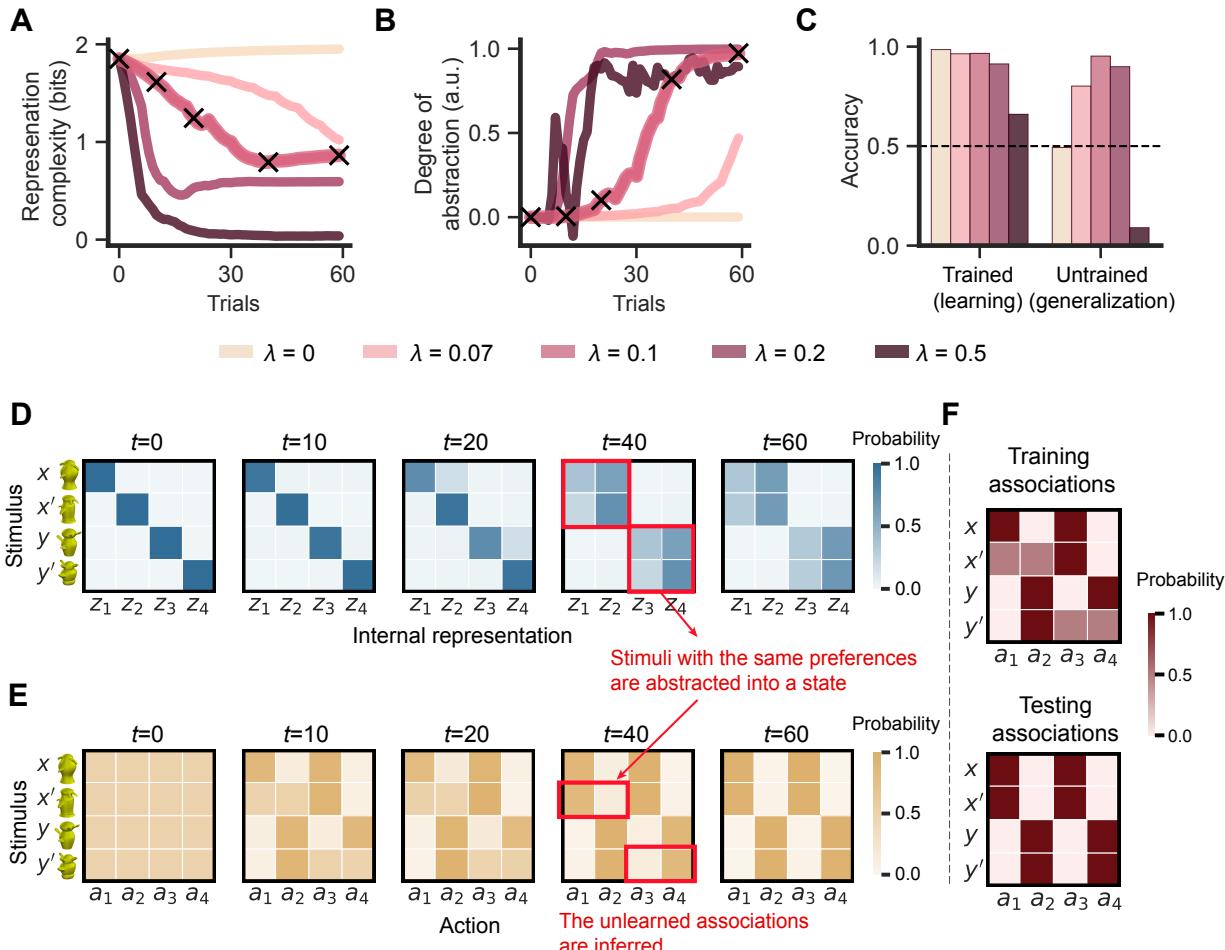
239 In the following sections, we will first demonstrate that, at the computational level, only
240 the ECPG model—which incorporates representation simplification—can qualitatively account
241 for human generalization behaviors (**Results 2.2 and 2.3**). We will then compare the ECPG model
242 to several published algorithmic-level models and show that, even without presuming any
243 specific algorithmic details about cognitive mechanisms, the ECPG model surpasses models with
244 handcrafted cognitive mechanisms in describing human behavior (**Result 2.4**). Overall, our
245 findings show that integrating efficient coding into the classical RL objective provides a more
246 comprehensive computational framework for understanding human learning and generalization.
247

248 **2.2 Abstract states inevitably merge in simplified representations, resulting in generalization**
249 Experiment 1 studies human generalization using the standard acquired equivalence paradigm.
250 In this setting, the four alien stimuli within each block are designed to have the same color but
251 different shapes and appendages (**Fig. 1D**). This design allows us to specifically study functional-
252 based generalization because the perceptual features (color, shape, and appendage) provide no
253 cues for generalization.

254 Why can humans generalize? The proposed efficient coding principle posits that, to
255 achieve simplified representations, an agent must appropriately abstract environmental stimuli
256 into robust latent states. Within each of these abstract states, the stimuli can then mutually
257 generalize. To illustrate this concept, we simulated the ECPG model at different levels of
258 simplicity, as determined by the parameter λ (0, 0.07, 0.1, 0.2, 0.5), while keeping other
259 parameters constant (See simulation details in **Method 4.10**). Note that when $\lambda = 0$, the ECPG
260 model collapses to the CAPG model, without efficient coding.

261 In these simulations, we first demonstrate that efficient coding forces state abstraction.
262 As shown in **Fig. 2A** ($\lambda = 0.1$), there is a significant decrease in the representation complexity of
263 the model from the beginning of training ($t = 0$) to the end ($t = 60$). This representation
264 simplification significantly affects the model’s internal representations ($\lambda = 0.1$). Before training,

265 when representations are complex, each stimulus is encoded in an unstructured way, with a one-
266 to-one correspondence in representation space (**Fig. 2D**, $t = 0$). Driven by efficient coding, the
267 ECPG model compresses representations, discarding redundant information and mapping stimuli
268 associated with the same actions into similar representations, resulting in the formation of
269 abstract states (**Fig. 2D**, $t > 20$, red arrows). We quantify the degree of the state abstraction using
270 the Silhouette score ([Rousseeuw, 1987](#)), which measures an object's (x) similarity to its own
271 latent state (x') relative to other states (y and y'). A score close to 0 indicates that stimuli are not
272 well-abstracted; a score close to 1 indicates that stimuli within each abstract state (x and x') are
273 encoded similarly and associate with each other, while stimuli across abstract states (x and y)
274 remain distinct. **Fig. 2B** ($\lambda = 0.1$) shows that the Silhouette score increases from 0, approaching
275 1, indicating that stable and meaningful abstract states automatically emerge from an initially
276 unstructured set of representations. The stimuli shared the same preferences became associated
277 with each other.



278

279 **Figure 2** The internal dynamics of the ECPG model during training. Panels A-E are generated
 280 by averaging over 400 simulations. Colors in panels A-C stands for different levels of simplicity,
 281 λ .

- 282 **A.** Representation complexity $I^\psi(S; Z)$ throughout the training process. The cross makers
 283 at $t = 0, 10, 20, 40, 60$ indicate the trials that are sampled for detailed analysis.
- 284 **B.** Throughout the training process, the effectiveness of state abstraction is measured by the
 285 silhouette score.
- 286 **C.** The proportion of correct responses for associations that were presented (trained) and
 287 not presented (untrained) during the training stage reflects learning and generalization
 288 performance, respectively. This figure includes only data from the testing stage. Dashed
 289 lines represent the 50% chance level.
- 290 **D.** Encoders $\psi(z|s)$ with $\lambda = 0.1$ at the sampled training trials. An encoder maps a stimulus
 291 $s \in \{x, x', y, y'\}$ to a distribution of internal representations $z_1 - z_4$. Each row stands for
 292 a categorical distribution that sums to 1. Darker shades indicate higher probability values.
 293 See **Extended Data 5.2** for encoders for other λ s.
- 294 **E.** Policies $\pi(a|s)$ with $\lambda = 0.1$ at the sampled training trials. A policy maps a stimulus $s \in$
 295 $\{x, x', y, y'\}$ to a distribution of actions $a_1 - a_4$. In each row, actions (a_1, a_2) and (a_3, a_4)
 296 form a probability distribution. See **Extended Data 5.2** for policies for other λ s.

297 **F.** Predefined correct associations in the training and testing stages. The dark tiles indicate
298 the correct associations, and the pink tiles stand for the associations not shown to the
299 participants

300

301 We then illustrate that stimuli within the same abstract states can generalize to each
302 other. After abstract states are stabilized ($t > 40$), the model begins to decode policy from the
303 structured representations, and changes in representation complexity become more nuanced
304 (**Fig. 2A**, $\lambda = 0.1$). Correspondingly, the policies decoded from stimuli within the same abstract
305 state are also similar (**Fig. 2E**, red arrows), illustrating that the ECPG model can effectively
306 generalize from training to testing associations. This is reflected by the model's significantly
307 above chance-level untrained accuracy (**Fig. 2C**, $\lambda = 0.1$), despite being exposed to only a subset
308 of the associations during training (**Fig. 2F**). Similar phenomena can be observed when λ is set to
309 0.2 (**Extended Data 5.3**).

310 Finally, we emphasize that the appropriate degree of state abstraction is critical; both
311 insufficient and excessive abstraction impair generalization performance. As λ increases, the
312 model shifts its focus from reward maximization to representation simplification, resulting in a
313 more intense and rapid state abstraction, yielding simpler states (**Fig. 2A**). When λ is set to lower
314 values ($\lambda = 0$ or 0.07), the ECPG model becomes more reward-focused and exhibits little or no
315 reduction in representation complexity (**Fig. 2A**, $\lambda = 0.07$ and 0). The insufficient compression
316 prevents the model from associating stimuli that share the same actions, leading to a failure in
317 state abstraction (**Fig. 2B**, $\lambda = 0.07$ and 0) and, consequently, compromised generalization
318 performance (**Fig. 2C**, $\lambda = 0.07$ and 0). Conversely, overly compressed representations (**Fig. 2A**, λ
319 $= 0.5$) tend to oversimplify abstraction, assigning all environmental states to a single internal state.
320 This abstraction can result in a significant reward loss, causing continuous adjustments to the
321 abstract states, as reflected by the oscillating Silhouette score (**Fig. 2B**, $\lambda = 0.5$). Such unstable
322 state abstraction can be detrimental to both generalization and learning performance (**Fig. 2C**, λ
323 $= 0.5$).

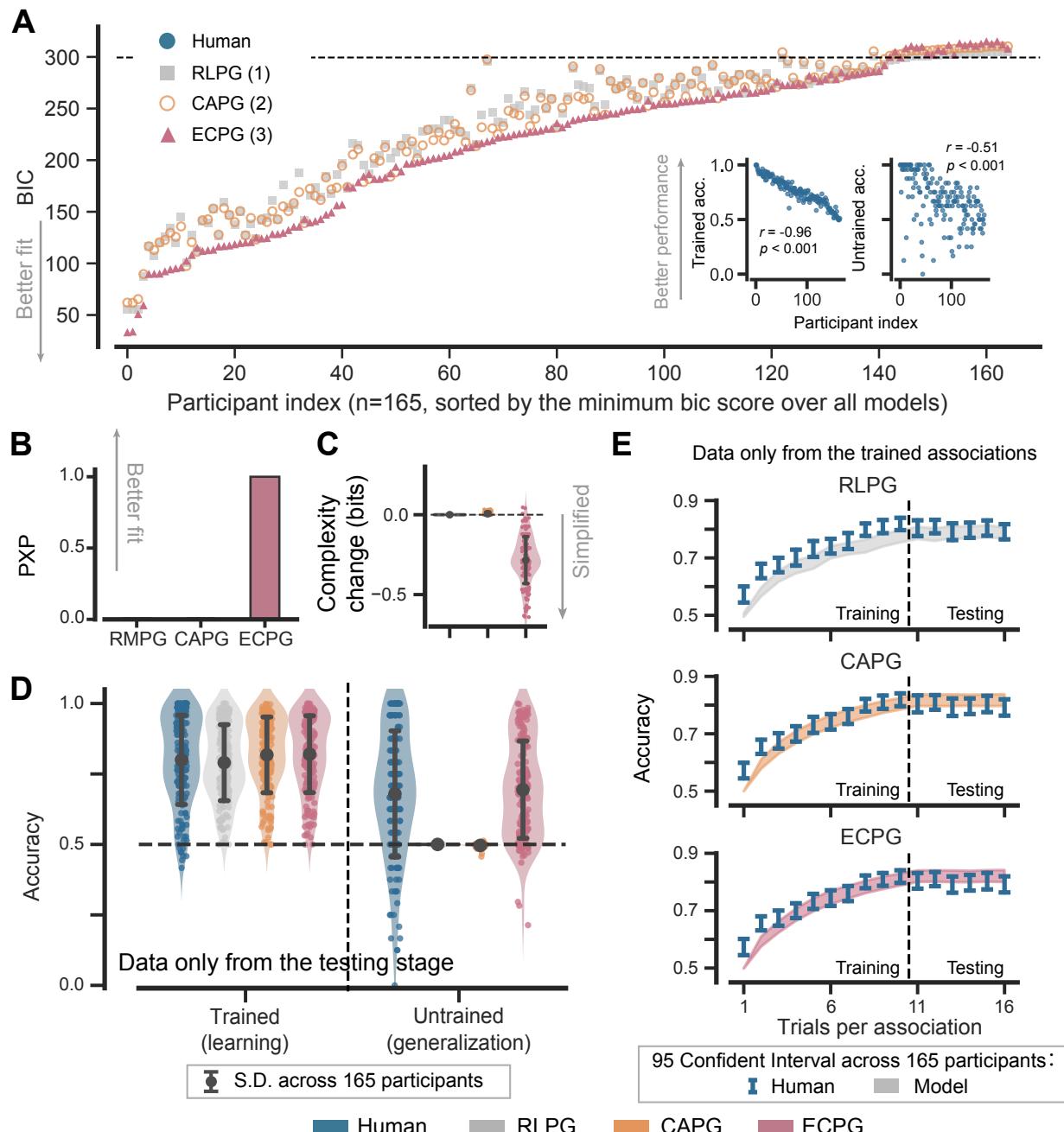
324 So far, our theoretical framework has outlined how efficient coding could result in
325 functional-based generalization. To verify whether these principles apply to the human brain, we
326 collected a set of behavioral data where 165 participants performed two blocks of standard
327 acquired equivalence tasks. We fitted all three models to the data and evaluated their

328 performance using the Bayesian Information Criterion (BIC). The ECPG model provided the best
329 description for a majority of participants (**Fig. 3A**). The model's fitting advantage was more
330 prominent in the testing stage, where generalization occurs, than in the learning stage (**Table 1**).
331 Some participants' behaviors could not be accurately described by the ECPG model, primarily due
332 to a lack of effort, which resulted in poor learning (Pearson's $r = -0.96, p < 0.001$) and
333 generalization performance (Pearson's $r = -0.51, p < 0.001$). We also conducted a Bayesian group-
334 level comparison and reported the protected exceedance probability (PXP)—the probability that
335 a model accounts for the data better than the others, beyond the chance level (the log model
336 evidence was estimated using the BIC) ([Rigoux et al., 2014](#)). As expected, the ECPG model again
337 was ranked first among the three models ($PXP > 0.999$; **Fig. 1B**). These findings underscore the
338 unique capability of the ECPG model in capturing human learning and generalization
339 performance.

340 To show that only the ECPG model learns simplified representations while the two control
341 models cannot, we computed the change in representation complexity during training. This was
342 quantified as the difference in the mutual information $I^\psi(S; Z)$ before and after training (**Fig.**
343 **3C**). The ECPG model successfully reduced the complexity of its representations, indicating that
344 it learned simplified representations as expected. In contrast, the two control models were
345 unable to compress their representations. Furthermore, we examined the simplicity parameter
346 λ of the ECPG model and found it to be significantly greater than 0 ($t(164) = 3.493, p < 0.001$,
347 Cohen's $d = 0.272$) (see **Extended Data 5.1** for model parameters). This finding suggests that the
348 representation simplification plays an important role in capturing human behaviors.

349 Human participants can effectively generalize in this task. Despite receiving no training,
350 the untrained accuracy for human participants is significantly greater than the 50% chance level,
351 though slightly lower than their accuracy for trained associations (**Fig. 3D**, blue). This observation,
352 consistent with numerous similar studies conducted previously ([Meeter et al., 2009](#); [Myers et al.,](#)
353 [2003](#); [Shohamy & Wagner, 2008](#)), indicates that human participants effectively generalized from
354 their prior learning. The ECPG model closely captures this generalization phenomenon, whereas
355 the two control models cannot generalize at all, with the untrained accuracy remaining at the 50%
356 chance level (**Fig. 3D**, gray and orange). More importantly, the strong performance of the ECPG

357 model in capturing human generalization performance does not come at the cost of its
 358 explanatory power in human instrumental learning behavior. It offers a description as precise as
 359 that of the two control models concerning the human learning curve throughout the training
 360 stage (**Fig. 3E**).



361
 362
 363

Figure 3 Behaviors of humans and models in Experiment 1. Models were fit to all behavioral data of each participant in both training and testing stages.

- 364 **A.** Models' Bayesian information criterion (BIC) for each participant. The RLPG model has 1
 365 parameter, the CAPG model has 2 parameters, and the ECPG model has 3 parameters.
 366 Also, see **Table 1** for the exact value.
- 367 **B.** Protected exceedance probability (PXP) tallies for each model.
- 368 **C.** Change in representation complexity after training. Scatterplot data located above the
 369 horizontal dashed line indicate representation expansion, while those below the line
 370 indicate representation compression. The RLPG model assumes that environmental
 371 stimuli are always perfectly reconstructed and, therefore, yield no change in complexity.
 372 Error bars reflect the standard deviation across 165 valid participants.
- 373 **D.** Proportion of correct responses for trained and untrained associations in the testing
 374 stage, representing learning and generalization performance, respectively. Only data
 375 from the testing stage is included. The dashed lines indicate the 50% chance level. Error
 376 bars reflect the standard deviation across 165 valid participants.
- 377 **E.** Proportion of correct responses over the number of trials for each association. The
 378 dashed line splits the experiment into two stages: on the left is the training stage, and on
 379 the right is the testing stage. Data from untrained associations is excluded from this
 380 learning curve analysis. Error bars reflect 95% confidence interval of the mean.
- 381

382 Based on both quantitative and qualitative evidence, we conclude that humans' ability to
 383 generalize originates from their computational goal of efficient coding. This process promotes
 384 the emergence of abstract latent states, which form the foundational basis for generalization.

385 **Table1. Model Fitting in Experiment 1**

	Overall			Training stage			Testing stage		
	RLPG(1)	CAPG(2)	ECPG(3)	RLPG(1)	CAPG(2)	ECPG(3)	RLPG(1)	CAPG(2)	ECPG(3)
NLL	114.562	111.16	101.202	66.196	64.235	62.061	48.366	46.925	39.141
AIC	231.124	226.319	208.403	134.392	132.47	130.121	98.732	97.85	84.282
BIC	234.499	233.07	218.529	137.179	138.045	138.484	101.52	103.425	92.644

386 **Bold values** indicate the best for each criterion

387 NLL=negative log likelihood; AIC=Akaike information criterion; Bayesian information criterion

388 The parameter number is shown in the brackets. For example, ECPG(3) means ECPG has 3 parameters.

389

390 **2.3 Efficient coding automatically extracts rewarding features throughout learning simplified
 391 representations**

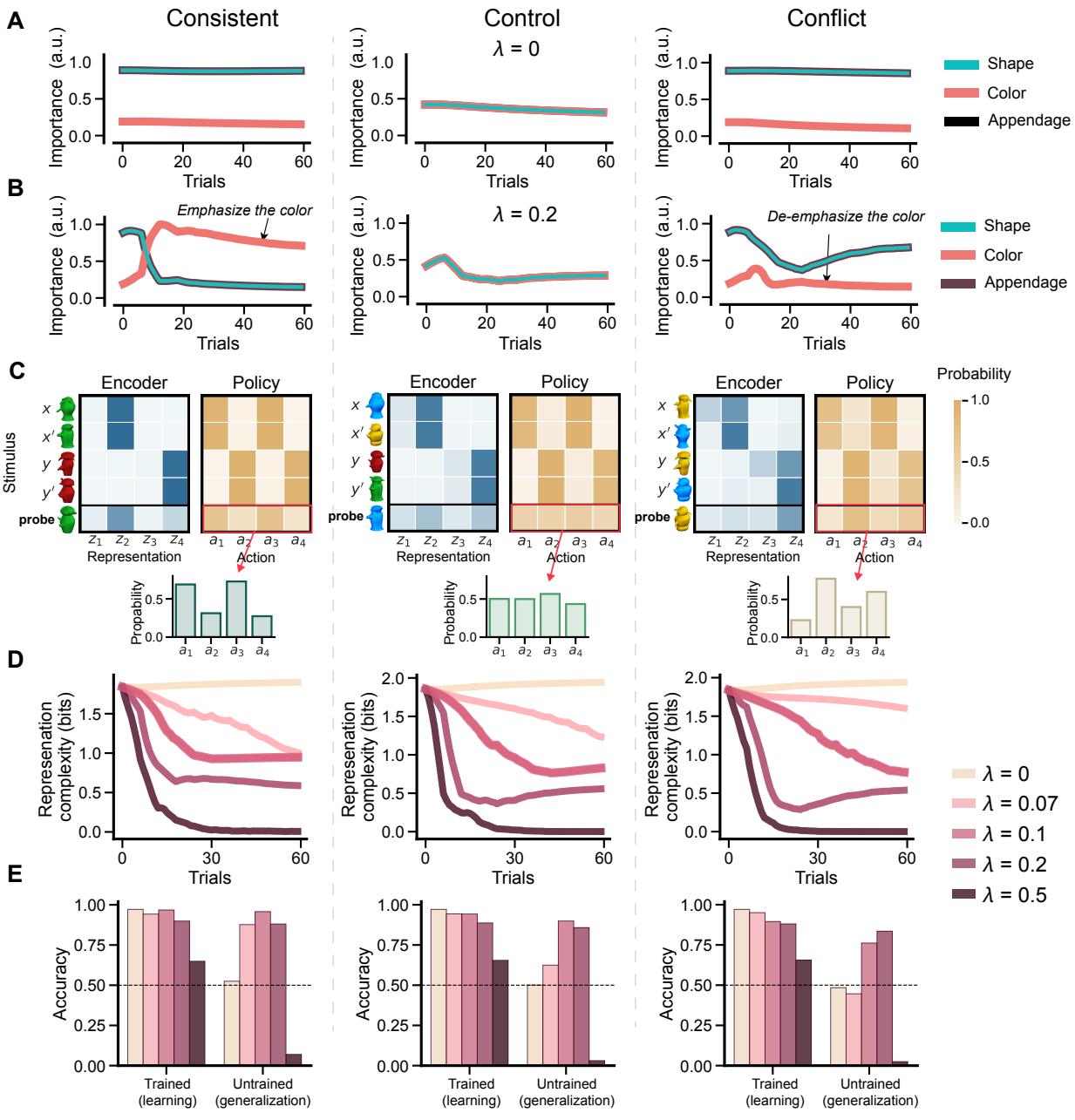
392 Experiment 2 extended the standard paradigm to examine both functional-based and perceptual-
 393 based generalizations of human participants. The experiment featured two primary
 394 modifications. First, we manipulated the stimuli's perceptual cues--shape, color, and appendage-
 395 -to ensure each feature provided a different amount of information about the environment's

396 rewards. We designed three experimental conditions, each with a distinct rewarding
397 configuration (**Fig. 1E**):

- 398 • In the *consistent* condition, the alien stimuli with the same color were associated with the
399 same actions, making the color the most rewarding feature.
- 400 • In the *control* condition, the colors of the stimuli were mutually different, and all features
401 were equally rewarding. This condition, like Experiment 1, only tested the functional-
402 based generalization and the state abstraction ability of an agent.
- 403 • In the *conflict* condition, stimuli with the same color were associated with different
404 actions, making shapes and appendages the rewarding features, while the color cue
405 yielded a negative reward.

406 These three conditions also indicated three levels of difficulty in rewarding feature extraction. In
407 the consistent and conflict condition, the four stimuli shared two colors, making color cues more
408 frequent and salient. For example, while the “cylinder” shape was associated with rewards twice
409 during training, the color “red” might have been rewarded four times. The consistent condition
410 was the easiest because this salient feature yielded positive rewards, while the conflict condition
411 was the most difficult because the agent needed to first suppress the color cue, the salient
412 feature, before being able to detect rewarding ones.

413 The second primary modification in Experiment 2 was the incorporation of a *probe*
414 stimulus during the testing stage; this stimulus was entirely new and had not been encountered
415 during training. This probe was used to assess humans’ ability to extract informative and
416 rewarding features at a behavioral level. A more detailed introduction to the use of this probe
417 design follows below, along with the presentation of our model’s predictions. All other aspects
418 of the experiment remained identical to those in Experiment 1.



419

420 **Figure 4** Predictions of the fECPG model in Experiment 2. All panels are generated by
 421 averaging over 400 simulations.

- 422 **A.** Simulated feature importance along with training of the fECPG model with $\lambda = 0$, which
 423 collapses to the fCAPG model. Note that the “shape” and “appendage” curves always
 424 overlap.
- 425 **B.** Simulated feature importance along with training of the fECPG model with $\lambda = 0.2$.
- 426 **C.** The predictive “probe” representation and policy of the fECPG agent ($\lambda = 0.2$) at the end
 427 of the training stage. An encoder maps a stimulus $s \in \{x, x', x, y'\}$ to a distribution of
 428 internal representations $z_1 - z_4$, and a policy maps a stimulus $ss \in \{x, x', x, y'\}$ to a
 429 distribution of internal representations $a_1 - a_4$. Darker tiles denote higher values. The
 430 policies for the control and conflict conditions are more stochastic than are those in the

431 consistent condition. The predictive policies applied to the probe stimuli are visualized in
432 both a heatmap and a bar plot.

433 **D.** Representation complexity throughout learning.

434 **E.** Learning and generalization performance for the fECPG model with different levels of
435 simplicity, $\lambda = 0, .07, 0.1, 0.2, 0.5$.

437 We reused the three models in Experiment 1, with only one key modification: adding a

438 *feature embedding function* to encode the stimuli's perceptual information. The feature
439 embedding function encoded each of the three visual features into a five-dimensional one-hot

440 code, where each dimension indicated a specific feature value. For example, the shape "cylinder"
441 was [1, 0, 0, 0, 0], the color "purple" was [1, 0, 0, 0, 0], and the color "yellow" was [0, 1, 0, 0, 0].

442 Each stimulus was represented by a combination of three such codes, concatenated into a 15-
443 dimensional vector to form the model's input. We refer to the models used in Experiment 2 as

444 feature RLPG (fRLPG; **Method 4.4.4**), feature CAPG (fCAPG; **Method 4.4.6**), and feature ECPG
445 (fECPG; **Method 4.4.5**) models to highlight their integration of the feature embedding construct.

446 We evaluated the model's feature extraction ability by analyzing the *importance* it

447 assigned to each feature, which quantifies each feature's contribution to the model's encoding
448 process ([Greydanus et al., 2018](#); [Guo et al., 2021](#)). Specifically, we added noise to one feature

449 dimension of the stimuli to perturb their appearance and examined the corresponding effect on
450 representations. A larger change in the representations indicated a higher importance of the
451 perturbed feature (see **Method 4.7**). Therefore, in this experiment, if a model consistently assigns
452 more importance to the predefined rewarding perceptual cue across all three experimental
453 conditions, we conclude that this model can effectively detect and extract rewarding features
454 within an environment.

455 Now that the stage is set, we can focus on answering three central research questions.

456 First, can the principle of efficient coding drive a model to extract rewarding features? Second, if
457 so, how do we validate that humans follow this principle in their learning processes? Third, how
458 does the rewarding feature extraction interact with the state abstraction ability examined in
459 Experiment 1?

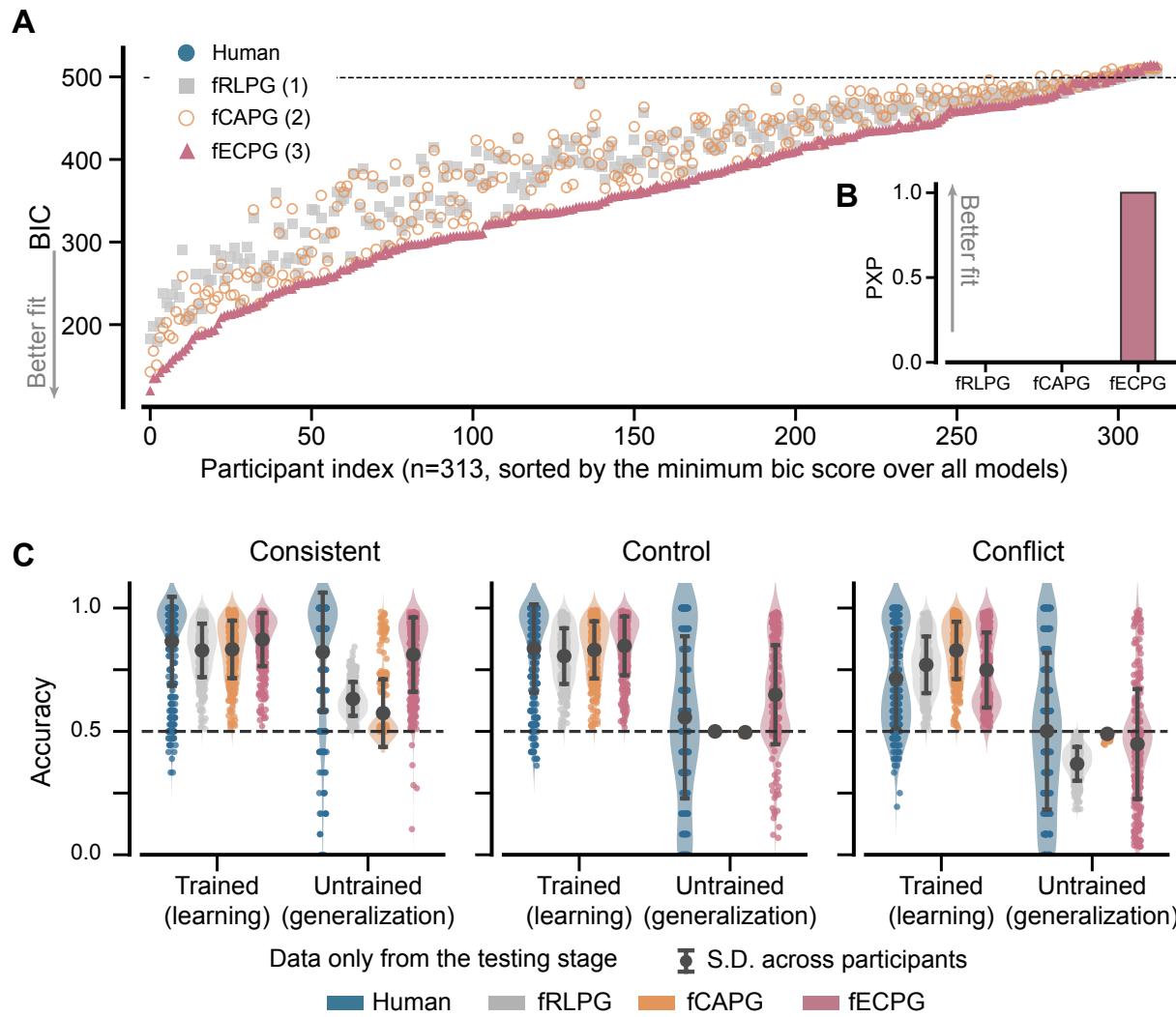
460 The answer to the first question is positive: efficient coding does promote rewarding

461 feature extraction. As designed in the experiment, color served as the rewarding feature in the
462 consistent condition, while it yielded negative rewards in the conflict condition. Driven by the

463 need for simpler representations—reflected in a focus on fewer features—the fECPG model ($\lambda =$
464 0.2) must selectively assign more importance to the color cue in the consistent case; and less
465 importance when the color becomes unrewarding (**Fig. 4B**, consistent). Conversely, in the conflict
466 condition, the model had to first deemphasize the salient cue color, due to its negative rewards,
467 and then reallocate importance to the other features contributing to positive rewards (**Fig. 4B**,
468 conflict). The demand for simplicity drives the model to focus on a subset of features, and the
469 goal of maximizing reward ensures that these focused features must be rewarding. In contrast, a
470 model without efficient coding ($\lambda = 0$) cannot adaptively reallocate feature importance during
471 its interaction with the environment. The model exhibits nearly the same feature importance
472 assignment for both consistent and conflict cases (**Fig. 4A**), indicating its inability to detect
473 rewarding information. It is worth noting that the fECPG model unintuitively predicts that shape
474 and appendage are rewarding features before training. We believe this is caused by our simplistic
475 approach to encoder initialization. We will further elaborate this point in the discussion section.
476 However, this observation does not undermine our conclusion.

477 To address the second question, we adopted a “probe” design within each block to
478 validate the theoretical predictions of the fECPG model on human participants. The probe
479 stimulus, introduced only during the testing stage and not present in the training, was designed
480 to always share the same color as stimulus x and the same shape as stimulus y' . According to
481 predictions of the fECPG model, humans should employ different policies in response to the
482 probe stimulus for different conditions. In the consistent condition, where color was the most
483 important feature, the probe stimulus should be perceived as similar to stimulus x (**Fig. 4C**,
484 consistent, encoder), leading to a response that coincides with the one for stimulus x (**Fig. 4C**,
485 consistent, policy). In this scenario, it is expected that human participants will demonstrate a
486 higher preference for actions a_1 and a_3 when responding to the probe stimuli. Conversely, in the
487 conflict condition where color was neglected, the probe stimulus should be perceived as more in
488 line with stimulus y' (**Fig. 4C**, conflict, encoder), which should be also reflected in the response
489 (**Fig. 4C**, conflict, policy). Therefore, human participants would be likely to use a_2 and a_4 . In the
490 control condition, given the lack of a dominant rewarding feature, the response to the probe

491 stimulus should not show a strong preference, being distributed between those for stimuli x and
 492 y' (Fig. 4C, control, policy).



507 For the third question, we observed that the efficiency of rewarding feature extraction
508 extends or shortens the time it takes to form stable abstract states, influencing the agent's
509 learning and generalization performance. In the consistent condition, the fECPG model quickly
510 identifies the rewarding feature, enabling it to form stable abstract states rapidly (**Fig. 4D**,
511 consistent), which results in a high degree of generalization (**Fig. 4E**, consistent). However, in the
512 control condition, where no salient cue dominates, the model experiences a slower state
513 abstraction process (**Fig. 4D**, control). In the conflict condition, the need to suppress color
514 prolongs the time required to extract rewarding features, thereby further extending the state
515 abstraction period (**Fig. 4D**, conflict). Consequently, the prolonged state abstraction period
516 reduces the time available for policy decoding, resulting in poorer learning and generalization
517 performance in the conflict case (**Fig. 4E**, control, conflict).

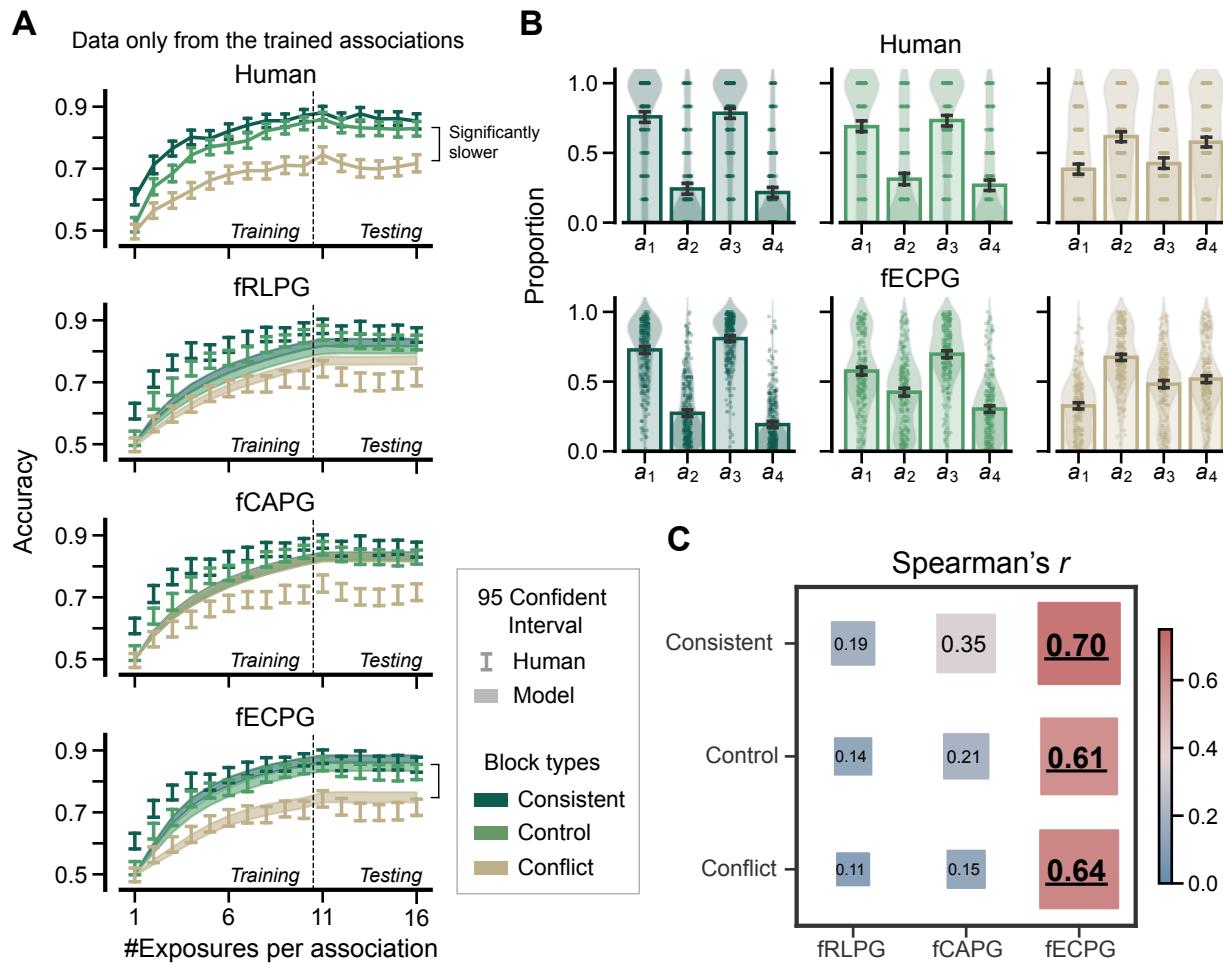
518 To validate these model predictions, we collected behavioral data from 313 participants
519 who each completed three task blocks corresponding to consistent, control, and conflict
520 conditions. We fit all three feature-based models and found that both BIC and PXP preferred the
521 fECPG model as the best model for capturing human behavioral data, consistent with the findings
522 from Experiment 1 (**Fig. 5A, B** and **Table 2**).

523 More importantly, as the principle of efficient coding predicts, human participants exhibit
524 different levels of generalization across experimental conditions. They achieved high untrained
525 accuracy in the consistent condition, lower in the control condition, and lowest in the conflict
526 condition (**Fig. 5C**, blue). Beyond the overall trend, humans' generalization behaviors are also
527 characterized by high variability. Some participants generalized effectively across all conditions,
528 while others always negatively transferred their knowledge. This variability is also accurately
529 captured by the fECPG model (**Fig. 5C**, red), but not by the two classical RL models without
530 efficient coding (**Fig. 5C**, gray and orange).

531 To our surprise, the ECPG model shows a significant advantage in predicting human
532 learning performance—an area where classical RL models have traditionally been preferred.
533 Participants demonstrated a more rapid improvement in the consistent condition than in the
534 control and conflict conditions. Specifically, the learning curve in the conflict condition was

535 markedly slower when compared to the other conditions. Only the fECPG model captured the
 536 significantly slower trend (**Fig. 6A**).

537 The probe design further validates the human participant's ability to extract rewarding
 538 features in a way predicted by the efficient coding principle. Human participants' responses to
 539 the probe stimuli were consistent with the fECPG model predictions (**Fig. 6B, C**; Spearman's $r >$
 540 0.60, $p < 0.001$ for all conditions; see **Method 4.9** for the correlation calculation). In contrast,
 541 models without efficient coding, fRLPG and fCAPG, failed to replicate such behavioral patterns
 542 (see **Extended Data 5.6** for their probe responses), exhibiting significantly weaker correlations
 543 with human behavioral data (**Fig. 6C**).



544

545 **Figure 6** Behaviors of human participants and models in Experiment 2 (part 2). Models were
 546 fit to all behavioral data of each participant in both training and testing stages.

547 **A.** Proportion of correct responses over the number of times that each association is shown.
 548 Dashed lines split the experiment into training and testing stages. The analysis excludes

549 responses to the untrained associations and the probe stimulus. Error bars indicate the
550 95% confidence interval for the mean estimate.

551 **B.** Humans and the fECPG model respond to the probe stimuli. Error bars indicate the 95%
552 confidence interval for the mean estimate.

553 **C.** Correlation between model predictions and human responses to probe stimuli. The
554 annotated value represents Spearman's correlation coefficient under different
555 experimental conditions.

556
557 It is important to note that there is a discrepancy between our prediction and human
558 behavior in the control condition. In response to the probe, human participants were likely to
559 use the policy of stimulus x rather than a random policy. This phenomenon could have arisen
560 from two potential factors. First, during training, the experiment might not have adequately
561 balanced the presentation frequency of the stimuli. Participants learned two associations with
562 stimulus x and one with stimulus y' (with the other association tested in the testing stage), which
563 implies that stimulus x was shown twice as frequently as was stimulus y' . Consequently,
564 participants might have adaptively adjusted their encoding and decision-making on these
565 statistics and placed more attention on stimulus x . Second, the color feature might have been
566 inherently more salient to humans. When the three features were equally informative,
567 participants may have naturally prioritized the color feature. However, this gap does not
568 undermine our conclusion that the fECPG model best captures human participants' responses to
569 the probe stimulus.

570 All evidence leads to one conclusion: during learning, humans strive to distill
571 representations into their simplest and most essential forms. Driven by this goal, human
572 participants learn representations using a small subset of rewarding features within their
573 environments. They further simplify these representations by abstracting them into compact,
574 lower-dimensional internal states, which naturally leads to generalization.

575

Table 2. Model Fitting in Experiment 2

		fRLPG (1)	fCAPG (2)	fECPG (3)	LC (3)	MA (3)	ACL (5)	fL1PG (3)	fL2PG (3)	fRNDPG (3)	fDCPG (3)
Overall	NLL	196.291	190.74	170.913	180.222	176.958	171.557	176.716	172.168	191.569	190.792
	AIC	394.583	385.48	347.826	366.445	361.917	353.114	359.433	350.336	389.138	387.585
	BIC	398.469	393.252	359.484	378.103	377.461	372.545	371.091	361.994	400.797	399.243
Consistent	NLL	56.465	55.154	47.357	51.416	48.993	44.073	48.151	47.412	55.628	55.166
	AIC	114.929	114.307	100.715	108.831	105.986	98.147	102.302	100.823	117.255	116.333
	BIC	117.717	119.882	109.077	117.194	117.136	112.084	110.664	109.186	125.618	124.695
Control	NLL	64.176	61.41	56.948	59.687	57.689	58.525	59.253	57.944	62.399	61.425
	AIC	130.351	126.821	119.897	125.375	123.379	127.05	124.506	121.887	130.798	128.851
	BIC	133.139	132.395	128.259	133.737	134.529	140.988	132.869	130.25	139.161	137.213
Conflict	NLL	75.651	74.176	66.607	69.119	70.276	68.959	69.312	66.813	73.542	74.201
	AIC	153.302	152.352	139.215	144.239	148.552	147.917	144.625	139.625	153.085	154.402
	BIC	156.089	157.927	147.577	152.601	159.702	161.855	152.987	147.988	161.447	162.764

577 **Bold values** indicate the best for each criterion

578 NLL=negative log likelihood; AIC=Akaike information criterion; Bayesian information criterion

579

580 **2.4 The human brain optimizes efficient coding to enhance learning and generalization**

581 A potential argument is that the classical RL objective is still a sufficient computational goal to
 582 explain human behavior once it provides the necessary cognitive mechanisms at the algorithmic
 583 level. We oppose this view for two reasons. First, all current algorithmic-level models with
 584 handcrafted cognitive mechanisms fail to capture human behaviors as effectively as the ECPG
 585 model. Second, these necessary cognitive mechanisms are inherently designed to simplify
 586 representations, essentially pursuing efficient coding.

587 We developed and compared three algorithmic-level models (**Fig. 7A**). The first model,
 588 the *Latent Cause* (LC) model ([Gershman et al., 2017](#); **Method 4.4.7**; see also [Collins & Frank, 2013](#);
 589 [Collins & Frank, 2016](#); [Gershman et al., 2010](#)) employs a hierarchical nonparametric Bayesian
 590 process to simulate human state abstraction ability. During the learning period, the LC model
 591 categorizes observed stimuli into latent clusters and learns the decision policy for these clusters.
 592 The second model, called the *Memory-Association* (MA) model (**Method 4.4.8**), memorizes all
 593 stimuli and their preferred actions, establishing associations between stimuli that share the same
 594 actions. These associations facilitate the inference of correct actions in untrained tasks, thereby

enabling generalization. The third model, *Attention at Choice and Learning* (ACL; ([Leong et al., 2017](#)); **Method 4.4.9**; see also ([Ballard et al., 2018](#); [Niv et al., 2015](#))), learns the value of each feature and calculates the feature importance based on these values. The model uses a linearly weighted feature value for decision-making. Notably, the LC and MA models emphasize state abstraction ability, whereas the ACL model is designed to extract and prioritize rewarding features. Both abilities could emerge by optimizing for the efficient coding goal, but in a different computational formulation.

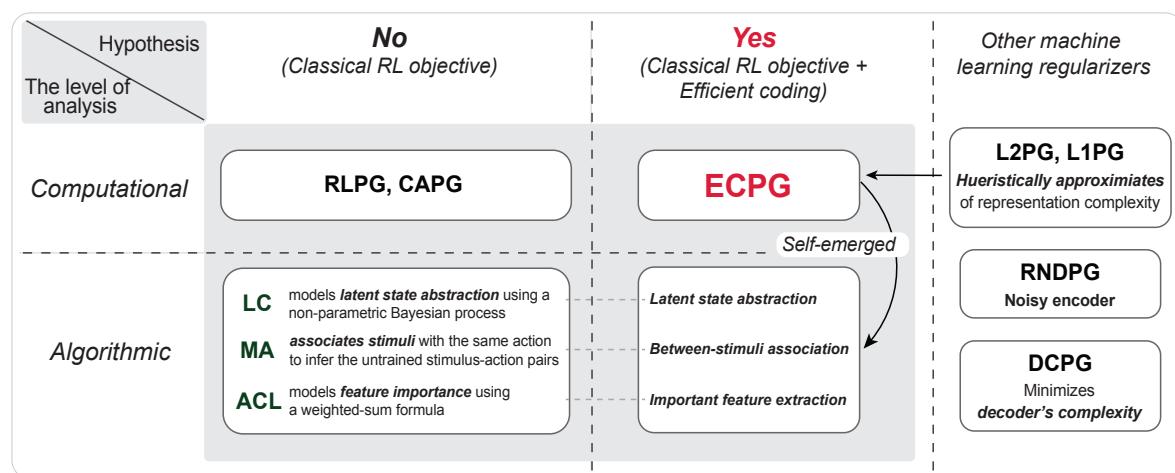
We tested all three algorithmic-level models on Experiment 2, with a focus on two qualitative metrics: generalization in the control case to examine their latent cause abstraction ability (**Fig. 7C**) and response to probe stimuli to evaluate their rewarding feature extraction (**Fig. 7D**). All three models underperformed the fECPG model in terms of BIC and PXP (**Fig. 7B, Table2**). The LC and MA models fail to account for human responses to probe stimuli (**Fig. 7D**) due to the lack of a feature extraction mechanism. The ACL model struggled with generalizing in the control case (**Fig. 7C**) as well as extracting rewarding features in both the control and conflict cases (**Fig. 7D**), because the feature importance calculated based on the accumulated value cannot deemphasize the negatively-rewarded feature effectively (see **Supplemental Note 2.2** for further discussion). These results underscore the superior performance of the fECPG model, a computational-level model, in modeling human behaviors and support our hypothesis that human participants learn simplified representations when maximizing rewards.

From a machine learning perspective, the fECPG model proposed here addresses a regularized optimization problem. This raises a final question: can the efficient-coding term be substituted by other commonly-used machine learning regularizers? We implemented an *L1-Norm Policy Gradient* (L1PG; **Method 4.4.10**) and an *L2-Norm Policy Gradient* (L2PG; **Method 4.4.10**), incorporating L1 or L2 norms as heuristic approximations for representation complexity. While the L1PG model exhibited a significant performance gap, the L2PG model performed more comparably to the fECPG model, as shown in **Fig. 7B, C, and D**. Although a substantial portion (~36%) of participants were better described by the L2PG model, these participants displayed distinct behavioral dynamics: compared to participants better captured by the fECPG model, they tended to learn more slowly and showed weaker generalization (see **Supplementary Note 2.3**

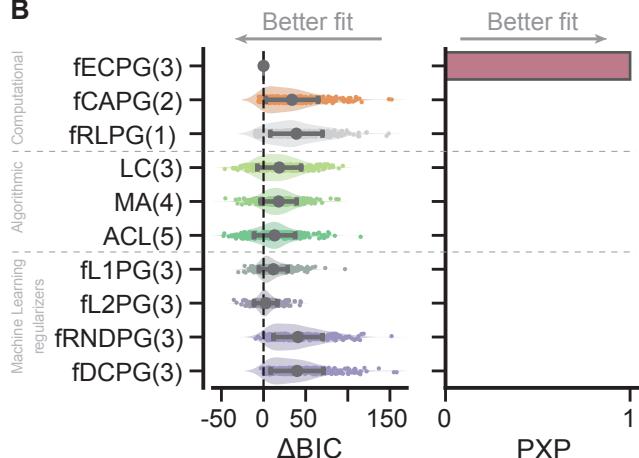
624 **for details**). This suggests that the ECPG model has a unique capability in capturing humans' fast
625 learning and strong generalization patterns. For completeness, we also tested a *Random*
626 *Regularizer Policy Gradient* (RNDPG; **Method 4.4.11**), which injects noise into the encoder
627 weights, as well as a Decoder Complexity Policy Gradient (DCPG; **Method 4.4.10**), which
628 constrains decoder complexity. However, these models proved ineffective at both generalizing
629 and extracting rewarding features (**Fig. 7B, C, D**).

630 Finally, we validated our conclusions by performing a model recovery analysis to test our
631 ability to differentiate between models (**Extended Data 5.8**). Importantly, we found that the
632 ECPG model can be uniquely distinguished from the other models. The low false positive rate
633 (with other models unlikely to be misidentified as ECPG) indicates that the ECPG model's superior
634 performance over the control models is not due to its expressiveness but to its accurate
635 description of human behavior. Thus, these findings confirm that our model recovery approach
636 robustly supports our conclusion that the ECPG model, with its augmented RL objective
637 incorporating efficient coding, best accounts for human learning and decision-making.

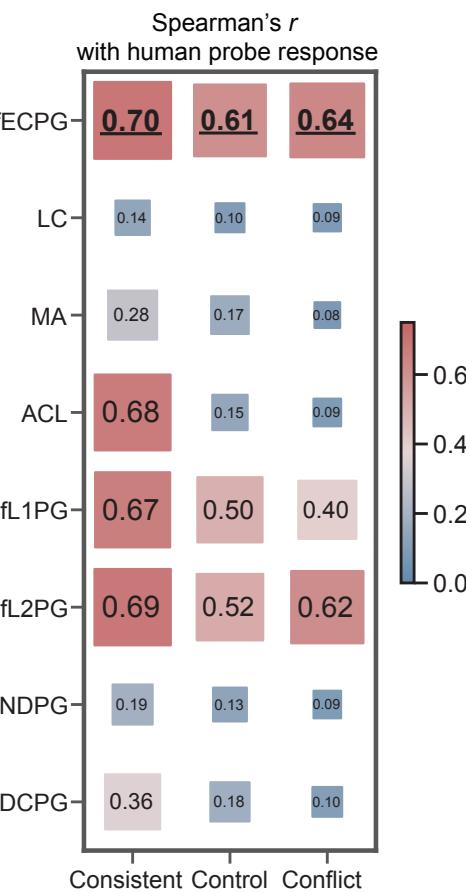
A Does the human brain optimize for efficient coding to achieve generalization?



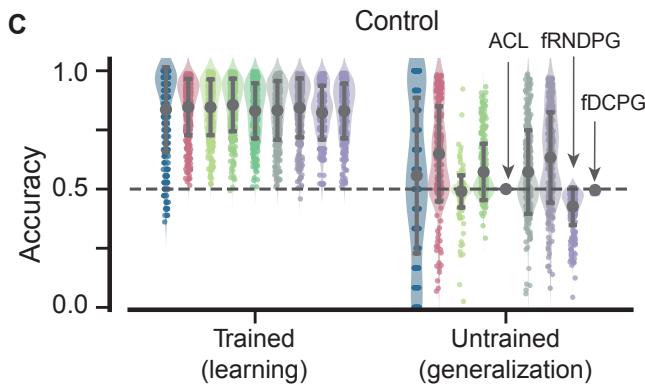
B



D



C



638

639 **Figure 7** Comparison of the (f)ECPG model with other algorithmic-level models in Experiment
 640 2. Error bars reflect the standard deviation (S.D.).

641 **A.** An overview of models across hierarchical levels. The central research question explores
 642 whether the human brain optimizes for efficient coding to enhance generalization. At the
 643 computational level, the ECPG model affirms this hypothesis ("Yes"), whereas the RLPG
 644 and CAPG models represent the opposing viewpoint ("No"). These models collectively
 645 represent the entire hypothesis space at the computational level. Below this, the ECPG

646 model is contrasted with various algorithmic-level models (LC, MA, ACL), each designed
647 with specific cognitive mechanisms. Additionally, the ECPG model is compared against
648 several common machine learning regularizers (L2PG, L1PG, DCPG) that also aim to
649 reduce model complexity, but through different methods.

- 650 **B.** Model comparisons for all models in terms of BIC and PXP in Experiment 2. See **Extended**
651 **Data 5.5** for model comparison in Experiment 1.
- 652 **C.** Learning and generalization performances across all models in the control case of
653 Experiment 2. See **Extended Data 5.6** for generalizations in other conditions.
- 654 **D.** Correlation between model predictions and human responses to probe stimuli at
655 different experimental conditions. See **Extended Data 5.7** for the bar plots.

656

657

658

659 **3. Discussion**

660 The classical RL framework has limitations in terms of its ability to explain human representation
661 learning and generalization. In this paper, we proposed augmenting the classical RL objective
662 with the efficient coding principle: an intelligent agent should distill the simplest necessary
663 representations that enable it to achieve its behavioral objectives. A computational-level model
664 derived from the revised framework (*efficient coding policy gradient*; ECPG), predicts that an
665 intelligent agent automatically learns to construct representations with a small set of rewarding
666 features with the environment. These representations are further simplified by abstracting them
667 into compact, lower-dimensional internal states, which naturally results in generalization. These
668 predictions were validated in two behavioral experiments, where the ECPG model consistently
669 provided a more accurate description of human behavior than two classical RL models without
670 efficient coding as well as several published human representation learning models. These
671 findings indicate that efficient coding offers a more suitable computational objective in
672 understanding human behavior.

673 In this paper, we examine whether the classical RL objective alone, or in combination with
674 efficient coding, better aligns with Marr's computational level in explaining human behavior. A
675 potential critique of our approach in **Section 2.4** is the lack of comparison with an alternative
676 model capable of generalizing without representation simplification. However, we found no such
677 model in the existing literature. This absence reflects the historical context of the acquired
678 equivalence paradigm on which our study builds. Although generalization within this paradigm
679 has long been documented since [Hall \(1991\)](#), previous explanations—including categorization
680 ([Urcuioli & Vasconcelos, 2008](#)), stimulus association ([Shohamy & Wagner, 2008](#)), and selective
681 attention ([Bonardi et al., 2005](#))—are all encompassed by our efficient coding framework. In other
682 words, algorithmic models based on selective attention, for example, inherently implement
683 mechanisms predicted by the computational-level goal of efficient coding.

684 The proposed ECPG model can potentially serve as a valuable quantitative analytical tool
685 for understanding cognitive impairments in individuals with mental disorders. Previous research
686 using the acquired equivalence paradigm has demonstrated that people with schizophrenia
687 ([Farkas et al., 2008](#); [Keri et al., 2005](#)), mild Alzheimer's disease ([Bodi et al., 2009](#)), hippocampal

688 atrophy ([Myers et al., 2003](#)), and Parkinson's disease ([Myers et al., 2003](#)) exhibit dysfunction in
689 the acquired equivalence task. We believe that the ECPG model, given its detailed portrayal of
690 human learning and generalization in this paradigm, can help elucidate the underlying causes of
691 these cognitive anomalies.

692 Beyond serving as a better empirical model for human learning, the proposed
693 computational objective could potentially represent a rational strategy (specifically, a resource-
694 rational strategy; see below) for humans. The classical RL objective was designed to maximize
695 expected reward in narrowly defined settings, where agents focus on learning a single, well-
696 defined task ([Sutton & Barto, 2018](#)). However, humans live in more complex and dynamic real-
697 world environments, where decision-making requires agents to generalize effectively from past
698 experiences to earn rewards in unseen scenarios. Moreover, the human brain is innately
699 capacity-constrained ([Miller, 1956](#)); it has inherent limitations in processing and storing
700 information, which requires the efficient use of cognitive resources. Therefore, learning simpler
701 representations that facilitate generalization is a crucial component in the pursuit of maximizing
702 reward in real-world decision-making. We believe that this insight can also improve learning and
703 generalization in artificial intelligence operating under real-world conditions.

704 The idea of linking RL to efficient coding has been applied to understand learning and
705 generalization in various contexts ([Berger & Machens, 2020](#); [Botvinick et al., 2015](#); [Franklin &](#)
706 [Frank, 2020](#); [Frydman & Jin, 2022](#); [Jaskir & Frank, 2023](#); [Luettgau et al., 2023](#); [Xia & Collins, 2021](#)).
707 For example, this approach has been shown to better explain monkeys' neural activity in frontal
708 areas ([Berger & Machens, 2020](#)), humans' risky choice behavior ([Frydman & Jin, 2022](#)), and meta-
709 level generalization between tasks ([Franklin & Frank, 2020](#)). Here, we present a specific
710 formalization of efficient coding using information-theoretic measures. We demonstrate that this
711 approach provides a better empirical description of both human learning and generalization
712 behaviors compared to several alternatives.

713 Our study bridges the gap between representation learning in the human brain and
714 machine learning. In cognitive science, researchers have applied latent cause clustering (LC) and
715 Association-Choice Learning (ACL) models to understand a variety of phenomena ([Radulescu et](#)
716 [al., 2021](#)). The latent cause clustering explains Pavlovian conditioning and extinction ([Gershman](#)

717 [et al., 2010](#)), memory modification ([Gershman et al., 2017](#)), social classification ([Gershman &](#)

718 [Cikara, 2023](#)), functional-based generalization ([Collins & Frank, 2013](#); [Collins & Frank, 2016](#);

719 [Lehnert et al., 2020](#)), etc. The selective attention, on the other hand, is used to explain concept

720 formation ([Mack et al., 2016](#)), the evolution of beliefs ([Markovic et al., 2015](#)), and has received

721 neural evidence from eye-tracking and functional Magnetic Resonance Imaging (fMRI) studies

722 ([Leong et al., 2017](#); [Niv et al., 2015](#)). In machine learning, researchers have focused on how

723 information-theoretic regularizers facilitate an artificial agent performing complex cognitive

724 tasks. For example, information-theoretic regularizers may help an agent learn robust state

725 abstractions that enhance learning speed ([Chelombiev et al.](#); [Islam et al., 2022](#); [Konidaris, 2019](#))

726 and form a simple but informative world model ([Ferns & Precup, 2014](#); [Rakelly et al., 2021](#)). Our

727 study demonstrates that in a simple cognitive task, both mechanisms serve the unified objective

728 of minimizing representation complexity, guided by an information-theoretic regularizer. This

729 finding facilitates communication between the two fields and contributes to a unified research

730 framework for understanding both machine and human intelligence. Building on this line of

731 thought, we plan to extend the current framework in future research to more complex task

732 settings, such as multi-step Markov Decision Processes (MDPs), and explore whether complex

733 human behaviors like planning and multi-task learning align with the predictions of information-

734 theoretic regularizers within machine learning.

735 Recent research has suggested that human intelligence is more accurately described by

736 the principle of *resource-rationality* ([Gershman et al., 2015](#); [Griffiths et al., 2015](#)) than by the

737 classical notion of *rationality* ([Von Neumann & Morgenstern, 1947](#)). The resource-rationality

738 principle emphasizes the need to consider computational costs in the pursuit of maximum

739 reward, building on the classical notion of rationality. The combination of efficient coding and

740 reward maximization principles applied in this study encapsulates the idea of resource-rationality,

741 with reward maximization representing the notion of rationality and representation complexity

742 representing computational costs. The basic idea is that information transmission in the brain

743 incurs significant metabolic costs, thus minimizing representation complexity (a quantification of

744 average information transmitted into the brain) serves as a reasonable proxy to minimize

745 computational costs ([Zenon et al., 2019](#)). Notably, while numerous studies have employed

746 resource-rationality to explain deviations from pure rationality in human behavior ([Gershman](#),
747 [2020](#); [Lieder & Griffiths, 2020](#); [Sims, 2016](#)), our research further emphasizes the advantages
748 conferred by the principle, particularly in accounting for state abstraction, rewarding feature
749 extraction, and generalization.

750

751 **Reference**

- 752 Asadi, A., Abbe, E., & Verdú, S. (2018). Chaining mutual information and tightening generalization
753 bounds. *Advances in Neural Information Processing Systems*, 31.
- 754 Ballard, I., Miller, E. M., Piantadosi, S. T., Goodman, N. D., & McClure, S. M. (2018). Beyond
755 Reward Prediction Errors: Human Striatum Updates Rule Values During Learning. *Cereb
756 Cortex*, 28(11), 3965-3975.
- 757 Barlow, H. B. (1961). Possible principles underlying the transformation of sensory messages.
758 *Sensory communication*, 1(01), 217-233.
- 759 Barto, A. G. (1995). Adaptive critics and the basal ganglia.
- 760 Bates, C. J., Lerch, R. A., Sims, C. R., & Jacobs, R. A. (2019). Adaptive allocation of human visual
761 working memory capacity during statistical and categorical learning. *J Vis*, 19(2), 11.
- 762 Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: a review and new
763 perspectives. *IEEE Trans Pattern Anal Mach Intell*, 35(8), 1798-1828.
- 764 Berger, S., & Machens, C. K. (2020). Compact task representations as a normative model for
765 higher-order brain activity. *Advances in Neural Information Processing Systems*, 33, 3209-
766 3219.
- 767 Bodí, N., Csibri, E., Myers, C. E., Gluck, M. A., & Keri, S. (2009). Associative learning, acquired
768 equivalence, and flexible generalization of knowledge in mild Alzheimer disease. *Cogn
769 Behav Neurol*, 22(2), 89-94.
- 770 Bonardi, C., Graham, S., Hall, G., & Mitchell, C. (2005). Acquired distinctiveness and equivalence
771 in human discrimination learning: evidence for an attentional process. *Psychon Bull Rev*,
772 12(1), 88-92.
- 773 Botvinick, M., Weinstein, A., Solway, A., & Barto, A. (2015). Reinforcement learning, efficient
774 coding, and the statistics of natural tasks. *Current opinion in behavioral sciences*, 5, 71-77.
- 775 Browning, M., Behrens, T. E., Jocham, G., O'Reilly, J. X., & Bishop, S. J. (2015). Anxious individuals
776 have difficulty learning the causal statistics of aversive environments. *Nat Neurosci*, 18(4),
777 590-596.
- 778 Chelombiev, I., Houghton, C., & O'Donnell, C. Adaptive estimators show information compression
779 in deep neural networks. arXiv 2019. *arXiv preprint arXiv:1902.09037*.
- 780 Collins, A. G., & Frank, M. J. (2013). Cognitive control over learning: creating, clustering, and
781 generalizing task-set structure. *Psychol Rev*, 120(1), 190-229.
- 782 Collins, A. G. E., & Frank, M. J. (2016). Neural signature of hierarchically structured expectations
783 predicts clustering and transfer of rule sets in reinforcement learning. *Cognition*, 152, 160-
784 169.
- 785 Crowston, K. (2012). Amazon mechanical turk: A research tool for organizations and information
786 systems scholars. Shaping the Future of ICT Research. Methods and Approaches: IFIP WG
787 8.2, Working Conference, Tampa, FL, USA, December 13-14, 2012. Proceedings,
- 788 Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-based influences
789 on humans' choices and striatal prediction errors. *Neuron*, 69(6), 1204-1215.
- 790 Farkas, M., Polgar, P., Kelemen, O., Rethelyi, J., Bitter, I., Myers, C. E., Gluck, M. A., & Keri, S.
791 (2008). Associative learning in deficit and nondeficit schizophrenia. *Neuroreport*, 19(1),
792 55-58.
- 793 Ferns, N., & Precup, D. (2014). Bisimulation Metrics are Optimal Value Functions. UAI,

- 794 Fisher, A., Rudin, C., & Dominici, F. (2019). All Models are Wrong, but Many are Useful: Learning
795 a Variable's Importance by Studying an Entire Class of Prediction Models Simultaneously.
796 *J. Mach. Learn. Res.*, 20(177), 1-81.
- 797 Flesch, T., Juechems, K., Dumbalska, T., Saxe, A., & Summerfield, C. (2022). Orthogonal
798 representations for robust context-dependent task performance in brains and neural
799 networks. *Neuron*, 110(24), 4212-4219.
- 800 Franklin, N. T., & Frank, M. J. (2020). Generalizing to generalize: Humans flexibly switch between
801 compositional and conjunctive structures during reinforcement learning. *PLoS Comput
802 Biol*, 16(4), e1007720.
- 803 Frydman, C., & Jin, L. J. (2022). Efficient coding and risky choice. *The Quarterly Journal of
804 Economics*, 137(1), 161-213.
- 805 Gershman, S. J. (2020). Origin of perseveration in the trade-off between reward and complexity.
806 *Cognition*, 204, 104394.
- 807 Gershman, S. J., Blei, D. M., & Niv, Y. (2010). Context, learning, and extinction. *Psychol Rev*, 117(1),
808 197-209.
- 809 Gershman, S. J., & Cikara, M. (2023). Structure learning principles of stereotype change. *Psychon
810 Bull Rev*, 30(4), 1273-1293.
- 811 Gershman, S. J., & Daw, N. D. (2017). Reinforcement learning and episodic memory in humans
812 and animals: an integrative framework. *Annual review of psychology*, 68, 101-128.
- 813 Gershman, S. J., Horvitz, E. J., & Tenenbaum, J. B. (2015). Computational rationality: A converging
814 paradigm for intelligence in brains, minds, and machines. *Science*, 349(6245), 273-278.
- 815 Gershman, S. J., Monfils, M. H., Norman, K. A., & Niv, Y. (2017). The computational nature of
816 memory modification. *Elife*, 6.
- 817 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- 818 Greydanus, S., Koul, A., Dodge, J., & Fern, A. (2018). Visualizing and understanding atari agents.
819 International conference on machine learning,
- 820 Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources: levels of
821 analysis between the computational and the algorithmic. *Top Cogn Sci*, 7(2), 217-229.
- 822 Guo, S. S., Zhang, R., Liu, B., Zhu, Y., Ballard, D., Hayhoe, M., & Stone, P. (2021). Machine versus
823 human attention in deep reinforcement learning tasks. *Advances in Neural Information
824 Processing Systems*, 34, 25370-25385.
- 825 Hall, G. (1991). *Perceptual and Associative Learning*. Oxford University Press.
- 826 Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., & Lerchner, A.
827 (2016). beta-vae: Learning basic visual concepts with a constrained variational framework.
828 International conference on learning representations,
- 829 Ho, M. K., Abel, D., Correa, C. G., Littman, M. L., Cohen, J. D., & Griffiths, T. L. (2022). People
830 construct simplified mental representations to plan. *Nature*, 606(7912), 129-136.
- 831 Huang, S., & Ontañón, S. (2020). A closer look at invalid action masking in policy gradient
832 algorithms. *arXiv preprint arXiv:2006.14171*.
- 833 Islam, R., Zang, H., Tomar, M., Didolkar, A., Islam, M. M., Arnob, S. Y., Iqbal, T., Li, X., Goyal, A., &
834 Heess, N. (2022). Representation learning in deep rl via discrete information bottleneck.
835 *arXiv preprint arXiv:2212.13835*.
- 836 Jaskir, A., & Frank, M. J. (2023). On the normative advantages of dopamine and striatal
837 opponency for learning and choice. *Elife*, 12.

- 838 Jiang, Y., Mi, Q., & Zhu, L. (2023). Neurocomputational mechanism of real-time distributed
839 learning on social networks. *Nat Neurosci*, 26(3), 506-516.
- 840 Keri, S., Nagy, O., Kelemen, O., Myers, C. E., & Gluck, M. A. (2005). Dissociation between medial
841 temporal lobe and basal ganglia memory systems in schizophrenia. *Schizophr Res*, 77(2-
842 3), 321-328.
- 843 Konidaris, G. (2019). On the necessity of abstraction. *Current opinion in behavioral sciences*, 29,
844 1-7.
- 845 Lehnert, L., Littman, M. L., & Frank, M. J. (2020). Reward-predictive representations generalize
846 across tasks in reinforcement learning. *PLoS Comput Biol*, 16(10), e1008317.
- 847 Leong, Y. C., Radulescu, A., Daniel, R., DeWoskin, V., & Niv, Y. (2017). Dynamic Interaction
848 between Reinforcement Learning and Attention in Multidimensional Environments.
849 *Neuron*, 93(2), 451-463.
- 850 Lerch, R., & Sims, C. R. (2021). Modeling associative motor learning through capacity-limited
851 reinforcement learning. *Journal of Vision*, 21(9), 2782-2782.
- 852 Li, F., Fergus, & Perona. (2003). A Bayesian approach to unsupervised one-shot learning of object
853 categories. *proceedings ninth IEEE international conference on computer vision*,
- 854 Li, L., Walsh, T. J., & Littman, M. L. (2006). Towards a unified theory of state abstraction for MDPs.
855 AI&M,
- 856 Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human cognition as
857 the optimal use of limited computational resources. *Behavioral and brain sciences*, 43, e1.
- 858 Lu, X., Lee, K., Abbeel, P., & Tiomkin, S. (2020). Dynamics generalization via information
859 bottleneck in deep reinforcement learning. *arXiv preprint arXiv:2008.00614*.
- 860 Luettgau, L., Erdmann, T., Veselic, S., Stachenfeld, K., Moran, R., Kurth-Nelson, Z., & Dolan, R.
861 (2023). Decomposing dynamical subprocesses for compositional generalization.
- 862 Mack, M. L., Love, B. C., & Preston, A. R. (2016). Dynamic updating of hippocampal object
863 representations reflects new conceptual knowledge. *Proceedings of the National
864 Academy of Sciences*, 113(46), 13203-13208.
- 865 Markovic, D., Glascher, J., Bossaerts, P., O'Doherty, J., & Kiebel, S. J. (2015). Modeling the
866 Evolution of Beliefs Using an Attentional Focus Mechanism. *PLoS Comput Biol*, 11(10),
867 e1004558.
- 868 Meeter, M., Shohamy, D., & Myers, C. E. (2009). Acquired equivalence changes stimulus
869 representations. *J Exp Anal Behav*, 91(1), 127-141.
- 870 Miller, G. A. (1956). The magical number seven plus or minus two: some limits on our capacity
871 for processing information. *Psychol Rev*, 63(2), 81-97.
- 872 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A.,
873 Riedmiller, M., Fidjeland, A. K., & Ostrovski, G. (2015). Human-level control through deep
874 reinforcement learning. *Nature*, 518(7540), 529-533.
- 875 Montague, P. R., Dayan, P., & Sejnowski, T. J. (1996). A framework for mesencephalic dopamine
876 systems based on predictive Hebbian learning. *J Neurosci*, 16(5), 1936-1947.
- 877 Myers, C. E., Shohamy, D., Gluck, M. A., Grossman, S., Onlaor, S., & Kapur, N. (2003). Dissociating
878 medial temporal and basal ganglia memory systems with a latent learning task.
879 *Neuropsychologia*, 41(14), 1919-1928.
- 880 Nelli, S., Braun, L., Dumbalska, T., Saxe, A., & Summerfield, C. (2023). Neural knowledge assembly
881 in humans and neural networks. *Neuron*, 111(9), 1504-1516 e1509.

- 882 Niv, Y. (2009). Reinforcement learning in the brain. *Journal of Mathematical Psychology*, 53(3),
883 139-154.
- 884 Niv, Y. (2019). Learning task-state representations. *Nat Neurosci*, 22(10), 1544-1553.
- 885 Niv, Y., Daniel, R., Geana, A., Gershman, S. J., Leong, Y. C., Radulescu, A., & Wilson, R. C. (2015).
886 Reinforcement learning in multidimensional environments relies on attention
887 mechanisms. *J Neurosci*, 35(21), 8145-8157.
- 888 Niv, Y., & Langdon, A. (2016). Reinforcement learning with Marr. *Curr Opin Behav Sci*, 11, 67-73.
- 889 Noh, H., You, T., Mun, J., & Han, B. (2017). Regularizing deep neural networks by noise: Its
890 interpretation and optimization. *Advances in Neural Information Processing Systems*, 30.
- 891 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,
892 N., & Antiga, L. (2019). Pytorch: An imperative style, high-performance deep learning
893 library. *Advances in Neural Information Processing Systems*, 32.
- 894 Pensia, A., Jog, V., & Loh, P.-L. (2018). Generalization error bounds for noisy, iterative algorithms.
895 2018 IEEE International Symposium on Information Theory (ISIT),
- 896 Rac-Lubashevsky, R., Cremer, A., Collins, A. G. E., Frank, M. J., & Schwabe, L. (2023). Neural Index
897 of Reinforcement Learning Predicts Improved Stimulus-Response Retention under High
898 Working Memory Load. *J Neurosci*, 43(17), 3131-3143.
- 899 Radulescu, A., Shin, Y. S., & Niv, Y. (2021). Human Representation Learning. *Annu Rev Neurosci*,
900 44, 253-273.
- 901 Rakelly, K., Gupta, A., Florensa, C., & Levine, S. (2021). Which Mutual-Information Representation
902 Learning Objectives are Sufficient for Control? *Advances in Neural Information Processing
903 Systems*, 34, 26345-26357.
- 904 Rescorla, R. A. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of
905 reinforcement and non-reinforcement. *Classical conditioning, Current research and
906 theory*, 2, 64-69.
- 907 Ribas-Fernandes, J. J., Solway, A., Diuk, C., McGuire, J. T., Barto, A. G., Niv, Y., & Botvinick, M. M.
908 (2011). A neural signature of hierarchical reinforcement learning. *Neuron*, 71(2), 370-379.
- 909 Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection for group
910 studies—revisited. *Neuroimage*, 84, 971-985.
- 911 Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster
912 analysis. *Journal of computational and applied mathematics*, 20, 53-65.
- 913 Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward.
914 *Science*, 275(5306), 1593-1599.
- 915 Shepard, R. N. (1987). Toward a universal law of generalization for psychological science. *Science*,
916 237(4820), 1317-1323.
- 917 Shohamy, D., & Wagner, A. D. (2008). Integrating memories in the human brain: hippocampal-
918 midbrain encoding of overlapping events. *Neuron*, 60(2), 378-389.
- 919 Schwartz-Ziv, R. (2022). Information flow in deep neural networks. *arXiv preprint
920 arXiv:2202.06749*.
- 921 Silver, D., Singh, S., Precup, D., & Sutton, R. S. (2021). Reward is enough. *Artificial Intelligence*,
922 299, 103535.
- 923 Simoncelli, E. P., & Olshausen, B. A. (2001). Natural image statistics and neural representation.
924 *Annu Rev Neurosci*, 24, 1193-1216.
- 925 Sims, C. R. (2016). Rate-distortion theory and human perception. *Cognition*, 152, 181-198.

- 926 Sims, C. R. (2018). Efficient coding explains the universal law of generalization in human
927 perception. *Science*, 360(6389), 652-656.
- 928 Sims, C. R., Jacobs, R. A., & Knill, D. C. (2012). An ideal observer analysis of visual working memory.
929 *Psychol Rev*, 119(4), 807-830.
- 930 Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- 931 Tishby, N., Pereira, F. C., & Bialek, W. (2000). The information bottleneck method. *arXiv preprint*
932 *physics/0004057*.
- 933 Tomov, M. S., Schulz, E., & Gershman, S. J. (2021). Multi-task reinforcement learning in humans.
934 *Nat Hum Behav*, 5(6), 764-773.
- 935 Urcuioli, P. J., & Vasconcelos, M. (2008). Effects of within-class differences in sample responding
936 on acquired sample equivalence. *J Exp Anal Behav*, 89(3), 341-358.
- 937 van Opheusden, B., Kuperwajs, I., Galbiati, G., Bnaya, Z., Li, Y., & Ma, W. J. (2023). Expertise
938 increases planning depth in human gameplay. *Nature*, 618(7967), 1000-1005.
- 939 Von Neumann, J., & Morgenstern, O. (1947). Theory of games and economic behavior, 2nd rev.
- 940 Wei, X. X., & Stocker, A. A. (2015). A Bayesian observer model constrained by efficient coding can
941 explain 'anti-Bayesian' percepts. *Nat Neurosci*, 18(10), 1509-1517.
- 942 Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist
943 reinforcement learning. *Machine learning*, 8, 229-256.
- 944 Xia, L., & Collins, A. G. E. (2021). Temporal and state abstractions for efficient learning, transfer,
945 and composition in humans. *Psychol Rev*, 128(4), 643-666.
- 946 Zenon, A., Solopchuk, O., & Pezzulo, G. (2019). An information-theoretic perspective on the costs
947 of cognition. *Neuropsychologia*, 123, 5-18.
- 948 Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., & Oliva, A. (2014). Learning deep 1143 features for
949 scene recognition using places database. *Z. Ghahramani*, 1144.
- 950
- 951

952 **4. Method**

953 **4.1 Ethics statement**

954 Participants were provided informed consent. The experimental protocol was approved by the
955 University Committee on Activities involving Human Subjects at Rensselaer Polytechnic Institute
956 (IRB-2055). Our experiment did not collect any demographical information from participants,
957 including gender.

958

959 **4.2 Materials**

960 The experiments were designed based on the paradigm of acquired equivalence (AE) ([Meeter et](#)
961 [al., 2009; Myers et al., 2003; Shohamy & Wagner, 2008](#)). The experiment included two types of
962 pictures: *alien* and *scene*. For the alien pictures, we utilized the “greebles” stimuli created by
963 Michael Tarr¹. The original greeble stimuli are purple. We created several new greeble stimuli by
964 changing their color. Regarding the scene pictures, we sampled from the “Places205” picture
965 database, as reported in Zhou et al., ([2014](#))².

966 In the main experiment blocks, aliens and scenes were organized into sets. Each set
967 comprised four alien stimuli (x, x', y, y') and four scenes (a_1, a_2, a_3, a_4), resulting in eight unique
968 associations per set. Six of these associations were trained during the association stage, while all
969 eight were tested in the testing stage. It is important to note that the stimuli in the AE task are
970 defined to be “superficially dissimilar”. In our experiment, the greeble stimuli within a block were
971 required to have the same color but exhibit mutually different shapes and appendages. There
972 was no correlation between the alien's shape (the configuration of its appendages) and the
973 correct response.

974

975 **4.3 Experiments**

976 **4.3.1 Experiment 1**

977 We recruited 302 participants from Amazon Mechanical Turk (MTurk) ([Crowston, 2012](#)). All
978 participants provided informed consent before starting the experiment. Each participant

¹ Greebles stimuli of Michael J Tarr, Carnegie Mellon University. Downloaded from <http://www.tarrlab.org/>

² Downloaded from <http://places.csail.mit.edu/downloadData.html>.

979 completed two practice blocks. To ensure a comprehensive understanding of the experiment,
980 participants were required to achieve at least 70% accuracy in the second practice block to
981 progress to the main experimental stage. Those who did not meet this criterion were allowed to
982 repeat the second practice block until they achieved the necessary performance level; otherwise,
983 they could not proceed to the main experiment. Participants received a base payment of \$2 plus
984 a bonus of up to \$3 based on their response accuracy.

985 This project aimed to study generalization within the learning process, meaning
986 participants who did not learn were outside the scope of this study. Consequently, we excluded
987 137 participants who failed a screening criterion (average accuracy lower than 60% for the last
988 24 trials, equating to 4 repetitions, in the training stage). All analyses in Experiment 1 were
989 conducted with the remaining 165 qualified participants.

990 The experiment consisted of two types of trials: training and testing. Each training trial
991 comprised three screens. Following a 500 ms fixation screen, the trial presented an alien stimulus
992 in the upper middle of the screen, along with photographs of scenes, offering one correct and
993 one incorrect choice. Participants were instructed, "Which scene is associated with this alien?"
994 and asked to respond by pressing the "F" or "J" key. These choices' left-right order was
995 counterbalanced across trials. The stimulus screen remained visible for ten seconds, followed by
996 a one-second feedback screen displaying either "Correct! You got 1 point." or "Incorrect! You got
997 0 point." The test stage trials were identical, except that no feedback was provided after
998 responses. In the testing stage, the experiment was directed to the next fixation screen following
999 the participant's response or after a maximum duration of ten seconds.

1000 Experiment 1 included two experimental blocks. Each block consisted of a training stage,
1001 during which participants learned the stimulus-action associations, and a testing stage, during
1002 which participants were tested on the learned associations as well as an untrained generalization
1003 probe (the dashed associations). The training stage involved each association being trained ten
1004 times with feedback, resulting in 6 (associations) \times 10 (repetitions) = 60 training trials. The testing
1005 stage tested both the trained and untrained associations six times, resulting in 8 (associations) \times
1006 6 (repetitions) = 48 testing trials. Participants were explicitly informed of the transition between

1007 the two experimental stages, and they were also reminded to keep and reapply their training
1008 experiences to achieve better performance.

1009 Before the main blocks, each participant was required to complete two practice blocks.
1010 The first practice block contained a simple trial-and-error learning task, where participants were
1011 trained to learn the correct answer through feedback. They were asked to correctly associate
1012 $x - a_1$ and $y - a_2$ without being asked to build any between-stimuli equivalence. This block
1013 provided a gentle introduction to the experiment, with ten trials and unlimited response time.
1014 The primary goal of the first practice was to familiarize participants with the trial-and-error
1015 training process. The second practice served as a quiz. This block included a simplified version of
1016 the main training stage, where participants were presented with four stimuli but only required
1017 to choose from two actions. It contained 4 (associations) \times 10 (repetitions) = 40 training trials.
1018 Participants needed to achieve 70% accuracy to pass the quiz; otherwise, they were required to
1019 repeat the second practice block before progressing to the main experimental blocks. The
1020 practice blocks were designed to help participants learn to establish between-stimuli equivalence
1021 in preparation for the main experimental blocks and used similar materials as the main blocks.

1022

1023 4.3.2 Experiment 2

1024 We recruited 497 participants from MTurk. All participants gave informed consent prior to the
1025 experiment. Each participant completed two practice blocks. To ensure full understanding of the
1026 experiment, they needed to achieve at least 70% accuracy in the second practice block to proceed
1027 to the main experimental stage. Those who did not meet this criterion were given the
1028 opportunity to repeat the second practice block until they reached the required accuracy. All
1029 participants received a \$3 base payment plus up to a \$4.5 bonus based on their response
1030 accuracy.

1031 We filtered the participants' data using the same screening criterion as in Experiment 1.
1032 A total of 184 participants were excluded because they did not achieve an average accuracy of
1033 60% for the last 24 trials (equivalent to 4 repetitions) in the training stage. All analyses in
1034 Experiment 2 were conducted with the remaining 313 qualified participants.

1035 Note that Experiment 2 included twice as many qualified participants as Experiment 1.
1036 This is because each participant in Experiment 1 completed two identical experimental blocks,
1037 while in Experiment 2, participants completed three different blocks, each corresponding to a
1038 different experimental condition. To ensure that each condition in Experiment 2 had a
1039 comparable amount of data to Experiment 1, we increased participant enrollment.

1040 After completing the same practice blocks as in Experiment 1, participants were required
1041 to complete three main experimental blocks: a consistent block, a control block, and a conflict
1042 block. The sequence of these blocks was counterbalanced among participants. The three blocks
1043 were almost identical; the only difference lay in the stimuli's appearance.

1044 Within each block, participants were required to complete a 60-trial training stage, which
1045 was the same as in Experiment 1. They then entered the testing stage, where they had to respond
1046 to eight regular testing associations plus an additional probe stimulus. Consequently, the testing
1047 stage comprised 9 (associations) \times 6 (repetitions) = 54 trials.

1048 The remaining details of Experiment 2 were identical to Experiment 1. Note that, unlike
1049 the untrained associations, we did not predefine a correct answer for the probe stimulus. We
1050 simply record participants' responses and hope to uncover which feature people were attending
1051 to by analyzing the response distribution.

1052

1053 **4.4 Models**

1054 To set the stage, we first formalize a dynamic decision process in the AE paradigm. For
1055 consistency, we adopt a notation system similar to that used in the experimental paradigm.

1056 We refer to a participant or decision maker as an *agent*. In each trial t , an agent is
1057 presented with an alien stimulus s_t from the set $\{x, x', y, y'\}$. The agent's task is to select an
1058 action, specifically a scene picture a_t , from the set $\{a_1, a_2, a_3, a_4\}$, with the objective of
1059 maximizing the reward $r(s_t, a_t)$ based on the feedback received. The subscript t denotes the
1060 variable at a particular trial. Both the stimulus S and action A are defined as categorical variables.

1061

1062 4.4.1 RLPG : Reinforcement learning policy gradient

1063 RLPG is a computational level model. The goal of the RLPG model is to identify a policy π that
 1064 optimizes the classical RL goal.

$$1065 \max_{\pi} E_{\pi} [r(s_t, a_t)] \quad (4)$$

1066 In the AE experiment, an agent was required to choose from two possible actions; before
 1067 receiving any feedback, each action had a 50% chance of being correct. The agent should have
 1068 had a baseline estimation of reward, denoted as b , prior to making a decision. An action is
 1069 considered positive when it yields a reward higher than the baseline and negative when the
 1070 reward is lower. We revised Eq. 4 to include this baseline reward estimation b ,

$$1071 \max_{\pi} E_{\pi} [r(s_t, a_t) - b] \quad (5)$$

1072 The formula indicates that the RL baseline learns to adjust the policy $\pi(a|s)$ to maximize received
 1073 reward subtracted by the baseline $r(s_t, a_t) - b$. The reward subtracted by the baseline is
 1074 commonly called *advantage* in the machine learning community. In this AE task, we assumed $b =$
 1075 0.5, corresponding to an expected reward of 0.5 (reward of 1 with 50% probability).

1076 The objective function can theoretically be tackled by any RL algorithm, but we have
 1077 chosen the simplest: the policy gradient method. We assume the policy follows a parameterized
 1078 softmax distribution, transforming the optimization problem into a parameter search:

$$1079 \pi(a|s; \phi) = \frac{\exp[\phi(s, a)]}{\sum_{a'} \exp[\phi(s, a')]} \quad (6)$$

1080 where ϕ denotes the parameters of the policy. Here, ϕ is a 4-by-4 table (4 stimuli by 4 action).
 1081 $\phi(s, a)$ refer the parameter a row s and column a (See **Extended Data 5.2** for model
 1082 architecture). For simplicity, we will denote this softmax formula as $\text{softmax}(\phi(s, a))$.

1083 Let $J(\phi) = \max_{\phi} E[r(s_t, a_t) - b]$, then the policy parameters were updated based on the
 1084 gradient of the objective function $\nabla_{\phi} J(\phi)$,

$$1085 \phi(s, a) = \phi(s, a) + \alpha_{\pi} \nabla_{\phi} J(\phi)(s, a) \quad (7)$$

1086 where $\alpha_{\pi} \geq 0$ is the learning rate of policy π . This policy learning rate is the only parameter in
 1087 the RL baseline model. The policy parameters ϕ were initialized to 0 before the experiment. Eq.
 1088 7 updates the policy via its gradient, which gives the name “policy gradient”. We have derived
 1089 the analytical gradient for both models and verified the derivation using pyTorch package

1090 (Paszke et al., 2019). See supplementary material for detailed derivation. The RLPG model
1091 features a single parameter: the policy’s learning rate, α_π .

1092 There are two remarks related to this simple model. First, though not explicitly shown,
1093 the RLPG assumes a perfect representation that fully reconstructs the stimulus. If we construct a
1094 model that explicitly includes the representation z and assume that each stimulus s
1095 deterministically maps to a unique representation z , the model nevertheless collapses to the
1096 RLPG model described above. Second, the RLPG model introduced in this study behaves similarly
1097 to the classic Q-learning model, extensively used in psychology (Daw et al., 2011). The most
1098 significant advantage of RLPG is its simplicity. The model has a single learning rate parameter,
1099 simultaneously approximating the effects of both the “learning rate” and “inverse temperature”
1100 parameters in the classic model. This allows for a more effective distillation of the computational
1101 essence underlying representation compression.

1102

1103 4.4.2 ECPG: efficient coding policy gradient

1104 The ECPG model is designed with a dual computational goal: to maximize reward while
1105 minimizing representation complexity.

$$1106 \max_{\psi, \rho} E_{\psi, \rho} [r(s_t, a_t) - b] - \lambda I^\psi(S; Z) \quad (8)$$

1107 The parameter $\lambda \geq 0$, referred to as the *simplicity parameter*, controls for the tradeoff
1108 between the classical RL objective and representation simplicity. When $\lambda = 0$, the agent does
1109 not compress stimuli representations for simplicity, focusing solely on reward maximization.
1110 Conversely, as $\lambda \rightarrow \infty$, the agent learns the simplest set of representations, encoding all stimuli
1111 into a single, identical representation. Therefore, an optimal λ balances compression and
1112 oversimplification.

1113 The introduction of latent representation z divides the policy into an encoder, ψ , and a
1114 decoder, ρ , both of which are optimized according to Eq. 8. Like the RLPG, we solve Eq. 8 using
1115 the policy gradient. Here, we considered a parameterized softmax encoder $\psi(z|s; \theta)$ and a
1116 decoder $\rho(a|z; \phi)$. The encoder parameter θ is a 4-by-4 table (4 stimuli by 4 representations) and
1117 the decoder parameter ϕ is also a 4-by-4 table (See **Extended Data 5.2** for model architecture).
1118 The policy π is derived from the combination of the encoder and decoder:

1119
$$\pi(a_t|s_t) = \sum_z \phi(z|s_t; \theta) \rho(a_t|z; \phi) \quad (9)$$

1120 We iteratively update the encoder and decoder to optimize Eq. 8 using the following
1121 scheme:

1122
$$\begin{cases} \max_{\theta} J(\theta) = \max_{\theta} \sum_z \psi(z|s_t; \theta) \left[\rho(a_t|z)(r(s_t, a_t) - b) - \lambda \log \frac{\psi(z|s_t; \theta)}{p(z)} \right] \\ \max_{\phi} J(\phi) = \max_{\phi} \sum_z \rho(a_t|z; \phi) [\psi(z|s_t)(r(s_t, a_t) - b)] \\ p(z) = \sum_s \psi(z|s)p(s) \end{cases} \quad (10)$$

1123 Here, $p(z)$ indicates the prior preference to the representation z . The first two optimization
1124 problems were solved using gradient ascent with learning rate parameters, α_ψ and α_ρ . The prior
1125 representation probability $p(z)$ was updated according to the definition of marginal probability.
1126 In practice, we also experimented with updating the prior in the gradient formula but found it
1127 made no significant difference in modeling human behavior. Therefore, we adopted the current
1128 scheme to reduce the number of free parameters.

1129 The initialization of the representation variable z is critical. In this article, z is a categorical
1130 variable that shares the same sample space as the stimulus. The encoder parameters θ are
1131 initialized by passing the product of an identity matrix and an initial value through a softmax
1132 function,

1133
$$\psi(z|s; \theta_0) = \text{softmax}(\theta_0 \mathbf{I}(s, z)) \quad (11)$$

1134 where $\mathbf{I}(s, z) = 1$ if $z = s$ otherwise 0. We pretrained the encoders to
1135 reach 99% discrimination accuracy by tuning the initial value θ_0 . The initial encoder in this case
1136 encoded stimuli almost orthogonally, reflecting the “superficially dissimilarity” in the definition
1137 of AE. See **Method 4.5** for more details.

1138 In summary, the ECPG model has three parameters: an encoder learning rate α_ψ , a
1139 decoder learning rate α_ρ , a simplicity parameter λ . The stimulus encoder parameters were
1140 initialized through pretraining, and the parameters of the decoder were initialized to 0.

1141 The ECPG model's encoder acts as a generative component, similar to the encoder in the
1142 beta variational autoencoder (β VAE) as described by Higgins et al., (2016), but with a categorical
1143 hidden layer instead of a continuous Gaussian distribution. This design facilitates the

1144 computation of mutual information and the quantification of representation complexity, building
1145 upon the work of Lu et al., ([2020](#)).

1146

1147 4.4.3 CAPG: the intermediary model

1148 The cascade policy gradient (CAPG) model is a special case of the ECPG model, with the simplicity
1149 fixed at 0, $\lambda = 0$. Therefore, the model has only two parameters: the learning rate for the
1150 encoder α_ψ and the learning rate for the decoder α_ρ .

1151

1152 4.4.4 fRLPG: the feature-based RLPG model

1153 The feature-based models we developed are extensions of the models previously introduced.
1154 The primary difference lies in the incorporation of a feature embedding function \mathcal{F} that maps a
1155 stimulus s onto a set of features f . We crafted a feature embedding function to decompose a
1156 greeble stimulus into three distinct features: shape, color, and appendage, using “one-hot
1157 encoding” for clear differentiation. For instance, the color “purple” is represented as [1, 0, 0, 0,
1158 0] and “yellow” as [0, 1, 0, 0, 0]. The feature function \mathcal{F} each input stimulus with its one-hot code
1159 and concatenates these codes into a 15-dimensional vector f , which serves as the model’s input
1160 (See **Extended Data 5.2C**).

1161 We modified the policy of the fRLPG model to create a feature-based baseline RL model
1162 that does not compress stimuli. This model proposes that visual similarity alone could account
1163 for human generalization performance, without the need for a representation compression
1164 mechanism. With the feature embedding function \mathcal{F} defined, the policy at trial t can be
1165 expressed as follows,

$$1166 \pi(a|\mathcal{F}(s_t);\phi) = \text{softmax}(\phi(f_t, a)) \quad (12)$$

1167 The parameter ϕ now is a 15-by-4 table. As before, the parameter ϕ was updated using
1168 the policy gradient method Eq. 7, and actions are selected by sampling from the softmax policy,
1169 Eq. 12. The model has one parameter, the learning rate of policy α_π .

1170

1171 4.4.5 fECPG: the feature-based ECPG model

1172 We modified the ECPG encoder to include the feature embedding function and preserved
 1173 the previous ECPG decoder formulation (See **Extended Data 5.2D**),

$$1174 \quad \psi(z|\mathcal{F}(s_t); \theta) = \text{softmax}(\theta(f_t, z)) \quad (13)$$

1175 The parameter θ now is a 15-by-4 table. The decoder $\rho(a|z; \phi)$ remains unchanged.

1176 A new challenge we faced was initializing the encoder parameters for the feature-based
 1177 model. The previous method, which relied on an identity matrix, was no longer suitable because
 1178 stimuli with overlapping features naturally appear more similar. For instance, a purple greeble
 1179 should be more similar to another purple greeble than to a yellow one.

1180 To address this, we introduced a new initialization technique:

1181 First, we measured the visual similarity between a stimulus s and all possible stimuli z
 1182 (including stimulus s itself) by calculating the dot product of their feature embeddings,

$$1183 \quad \text{sim}(s, z) = \langle \mathcal{F}(s), \mathcal{F}(z) \rangle \quad \forall z \in \{x, x', y, y'\} \quad (14)$$

1184 Second, we multiplied these similarity scores by a scalar θ_0 and passed them through a
 1185 softmax function to form the representation of stimulus s ,

$$1186 \quad \bar{\psi}(z|s) = \text{softmax}(\theta_0 \text{sim}(s, z)) \quad (15)$$

1187 As before, the initial value θ_0 is tuned through pretraining. This value controls the
 1188 perceived similarity between stimuli. When θ_0 is small, stimuli with overlapping features look
 1189 similar.

1190 Finally, we used the representation $\bar{\psi}(z|s)$ as supervised labels to train the model
 1191 encoder $\psi(z|\mathcal{F}(s), \theta)$ by minimizing their cross-entropy loss

$$1192 \quad \min_{\theta} - \sum_z \bar{\psi}(z|s) \log \psi(z|\mathcal{F}(s), \theta) \quad (16)$$

1193 Once the training has converged, the encoder for the fECPG model is considered
 1194 initialized. The subsequent learning and decision-making processes are consistent with the
 1195 original ECPG model.

1196 The model has three parameters $\{\alpha_{\psi}, \alpha_{\rho}, \lambda\}$: the learning rate for the encoder α_{ψ} , the
 1197 learning rate decoder α_{ρ} , and the simplicity parameter λ .

1198

1199 4.4.6 fCAPG: the feature-based CAPG model

1200 The fCAPG model is special case of the fECPG with the simplicity parameter fix to 0, $\lambda = 0$. The
1201 model has two parameters α_ψ and α_ρ : the learning rate for the encoder α_ψ , the learning rate for
1202 decoder α_ρ .

1203

1204 4.4.7 LC: latent-cause clustering

1205 The LC model is an algorithmic-level model, adopted and modified from Gershman et al., (2017).
1206 The central idea of the LC model is to use a non-parametric Bayesian process to model the
1207 cognitive process of latent-cause clustering. The original model cannot be directly applied to the
1208 acquired equivalence task, as it is a model about association learning, not instrumental learning.
1209 We modified it to incorporate a value function of latent-cluster and action $Q(z, a)$. To facilitate
1210 derivation, we rewrite this value function in the probability format $p(q|z, a)$.

1211 The heart of the LC model lies in learning a Chinese restaurant process to model the
1212 human latent-cause clustering process. To construct the LC model, we rely on the following five
1213 constructs:

- 1214 • τ : a weight parameter that decides the influence of a variable at time l to the current
1215 variable t , the weight follows a power function,

$$1216 \quad \tau_l = \frac{1}{t-l}, \quad \forall l \leq t-1 \quad (17)$$

- 1217 • $p(z|\mathbf{z}_{1:t-1})$: the Chinese restaurant prior of the latent-cause z :

$$1218 \quad p(z=k|\mathbf{z}_{1:t-1}) = \begin{cases} \frac{\sum_{l < t-1} \tau_l \mathbf{I}(z_l = k)}{\sum_{l < t-1} \tau_l + \alpha}, & \text{if } k \text{ is a old cause} \\ \frac{\alpha}{\sum_{l < t-1} \tau_l + \alpha}, & \text{if } k \text{ is a new cause} \end{cases} \quad (18)$$

1219 where $\mathbf{I}(\cdot)$ is an indicator function that returns 1 if the condition is satisfied otherwise 0.
1220 α is the concentration parameter of the Chinese restaurant process. In the following
1221 derivation, we will also use $p(z)$ to represent this component for short.

- 1222 • $p(s_t|z)$: the likelihood of the current stimulus given history:

$$1223 \quad p(s_t|z) = \sum_f \mathbf{I}(f = \mathcal{F}(s)) p(f|z) \quad (19)$$

The feature f has d dimension ($d = 3$ in this study stands for shape, color, appendage). Each dimension represents an index j of the feature (e.g. $f_2 = 1$ means yellow color). We slightly abused the notation $\mathcal{F}(s_t)$. The probability mass function of a feature vector f is written as,

$$p(f_1, \dots, f_d | z = k) = \frac{\prod_d \prod_j \lambda_{kd}^{\mathbf{I}(f_d=j)} + p}{\sum_d \prod_d \prod_j \lambda_{kd}^{\mathbf{I}(f_d=j)} + p} \quad (20)$$

where d is the dimension of feature vector f . p is the prior for the multi-categorical distribution and is always set to 0.1, $p = 0.1$, in this study. This likelihood distribution is learned throughout training by,

$$\lambda_{kd} = \frac{\sum_{l < t-1} \mathbf{I}(f_d = j)}{\sum_{d', i'} \sum_{l < t-1} \mathbf{I}(f_{d'} = j)} \quad (21)$$

- $p(q|z, a_t)$: the value of latent-cluster and action

$$p(q|z, a_t) = \mathbf{I}(q = \mathbf{w}_{z,a_t}) = Q(z, a_t) \quad (22)$$

The weight w is always initialized as $1/|A|$.

- $p(r_t|q)$: the reward distribution given q value:

$$p(r_t|q) = q^r(1-q)^{1-r} \quad (23)$$

This is a Bernoulli distribution, given that the experiment yields a deterministic reward of 0 or 1.

The LC model's learning procedure can be framed as a structural learning problem, which is addressed by the Expected-Maximization (EM) algorithm. See **Supplemental Note 1** for details.

- **E-step:** $p(z|s_t, a_t, r_t; \mathbf{w}) \propto \sum_q p(z)p(s_t|z)p(q|z, a_t)p(r_t|q)$
 - **M-step:** $\mathbf{w} = \mathbf{w} + \eta p(z|s_t, a_t, r_t; \mathbf{w})(r - Q(z, a_t))\mathcal{F}(s_t)$

where η is the learning rate.

The LC model needs to make two decisions, assigning the current stimulus s_t to a latent-cluster z_t and picking an action a_t for current stimulus s_t . The latent-cluster z_t is selected to maximize the posterior of the latent-cluster,

$$z_t = \arg \max_z p(z|s_t, a_t, r_t; \mathbf{w}) \quad (24)$$

1251 The stimulus value is passed through a softmax to generate the action. The stimulus value is
1252 determined by marginalizing all possible latent causes.

1253

$$\pi(a|s_t) = \text{softmax} \left(\beta \sum_z p(z)p(s_t|z)Q(z, a_t) \right) \quad (25)$$

1254 where β is the inverse temperature.

1255 The four parameters $\{\alpha, \eta, \beta, w_0\}$ of the LC are: the concentration parameter for Chinese
1256 restaurant process α , the learning rate η , the inverse temperature of the softmax β , and the
1257 initial weight for the value function w_0 .

1258

1259 **4.4.8 MA: memory-association model**

1260 The MA model is an algorithmic-level model that combines memory and association mechanisms.
1261 It memorizes stimuli-action pairs and forms associations between stimuli with shared actions or
1262 features, using these associations to infer actions for untrained tasks. Stimuli with salient shared
1263 features (like “color”) are more easily associated.

1264 During training, the MA model utilizes a feature-action Q value table to store and update
1265 the value of stimuli-action pairs based on prediction errors (similar to classical RL models). The
1266 update formula is as follows:

1267

$$Q(f, a_t) = Q(f, a_t) + \alpha_Q \delta_t, \quad \forall f \in \mathcal{F}(s) \quad (26)$$

1268 where the value stimulus action pair is calculated using $Q_s(s, a) = \sum_{f \in \mathcal{F}(s)} Q(f, a_t)$.

1269 The model also maintains an association table $w(s = i, s = j)$, that records the
1270 associations between pairs of stimuli based on shared actions. The update rule for this
1271 association strength is:

1272

$$w(i, j) = w(i, j) + \alpha_{\text{assoc}} (1 + k(i, j)) (\bar{w}(i, j) - w(i, j)) \quad (27)$$

1273 where the update target $\bar{w}(i, j) = \sum_{T(s, a')=1} Q_s(i, a')Q_s(j, a')$ evaluates the similarity of the
1274 action preferences between i and j within the trained association, and $T(s, a)$ is another table
1275 that memories whether a certain association $s - a$ has been trained before. k here measures the
1276 perceptual similarity between i and j using the dot product of their feature vectors, $k(i, j) =$
1277 $\langle \mathcal{F}(i), \mathcal{F}(j) \rangle$. α_{assoc} is the learning rate for the stimuli association matrix. Note that we weighted

1278 the association update using $1 + k(i, j)$ rather than $k(i, j)$ to ensure that the association update
1279 is not zero. The self-association is always 1, $k(i, i) = 1$.

1280 In the decision-making stage, the MA model first needs to decide whether the presented
1281 stimulus-action pair (s_t, a) has been trained before. For the trained association pairs, the model
1282 uses the learned state-action values for decision-making, while for the untrained pairs, the model
1283 infers their values based on the state-action values for all stimuli weighted by their associations
1284 $w(s_t, s)$,

$$1289 Q_s(s_t, a) = \begin{cases} \sum_{f \in \mathcal{F}(s_t)} Q(f, a_t), & \text{if } T(s_t, a) = 1 \\ \sum_s \phi(s|s_t, w) \sum_{f \in \mathcal{F}(s_t)} Q(f, a_t), & \text{otherwise} \end{cases} \quad (28)$$

1285 where $\phi(s|s_t, w) = \text{softmax}(\beta_{\text{assoc}} w(s_t, s))$ calculates the weights for each stimuli s , including
1286 s_t itself. The β_{assoc} is the inverse temperature. Stimuli with a higher association gain more weight
1287 in estimating the state-action value for the current stimulus s_t . $T(s_t, a) = 1$ means the stimulus-
1288 action association has been trained before.

1290 Like many other RL models (such as LC, ACL in this study), the MA model creates a policy
1291 by passing the state-action values through a softmax function,

$$1292 \pi(a|s_t) = \text{softmax}(\beta Q_s(s_t, a)) \quad (29)$$

1293 In total, the MA model has 4 parameters: the learning rate for the feature-action value
1294 function α_Q , the learning rate for the stimuli association matrix α_{assoc} , the inverse temperature
1295 for the association-based weights β_{assoc} , and the inverse temperature for the policy β .

1296

1297 4.4.9 ACL: attention at choice and learning

1298 The ACL model is an algorithmic-level model that model humans' rewarding feature extraction
1299 ability using a linear selective attention mechanism. The model was from Leong et al., (2017)
1300 with two modifications. First, the original ACL model was developed on a different paradigm and
1301 could not be directly applied to the current generalization task. We modified the ACL model to
1302 include a feature-action value $Q(f, a)$ design as Ballard et al., (2018). Second, instead of using
1303 the attention weights calculated from the eye-tracking and functional MRI data, we estimated
1304 the attention weight using an attention model. The original authors constructed two types of

1305 models to examine the bidirectional relationship between learning and attention. The “choice
 1306 models” utilize attention data, collected through eye tracking and fMRI, to predict human
 1307 behaviors. Conversely, the “attention models” use human behavioral data as input and predict
 1308 the recorded attention data. In this study, we implemented the ACL model by combining the best
 1309 choice model (the ACL model in the original paper) and the best attention model (the VALUE
 1310 model in the original paper). We chose this approach because our study lacks attention data,
 1311 such that we have to use the best attention model to provide reasonable estimation of the
 1312 attention weights.

1313 The action value of a stimulus s_t is a weighted sum of the feature-action value (Equation
 1314 1 in original paper),

$$Q(s_t, a_t) = \sum_{f \in \mathcal{F}(s_t)} \phi(f) Q(f, a_t) \quad (30)$$

1315 where ϕ indicates the selective attention on each feature. The estimation of ϕ will be shown
 1316 soon in Eq. 34. Note that, here we slightly abused the notation $\mathcal{F}(s_t)$. This feature embedding
 1317 function in the ACL model is no longer returning a feature vector. Instead, it returns all features
 1318 contains in a stimulus s . We did this to better align with the equations written in Leong et al.,
 1319 (2017). When receiving a reward r_t , the prediction error δ_t can be calculated as (Equation 2 in
 1320 original paper),

$$\delta = r_t - Q(s_t, a_t) \quad (31)$$

1321 which is then used to update the value of the chosen feature and action (Equation 3 in original
 1322 paper),

$$Q(f, a_t) = Q(f, a_t) + \eta \phi(f) \delta_t, \forall f \in \mathcal{F}(s_t) \quad (32)$$

1323 where η is the learning rate.

1324 The ACL model generates actions by passing the stimulus-action value through a softmax
 1325 function (Equation 4 in original paper).

$$\pi(a|s_t) = \text{softmax}(\beta Q(s_t, a)) \quad (33)$$

1326 where β is the inverse temperature.

1327 To estimate the attention, the original paper used the VALUE model. This model presumes
 1328 that attention tracks feature values, $\tilde{Q}(s, a)$. Note that this $\tilde{Q}(s, a)$ is a new value function used
 1329 to estimate attention weight, independent from the Q function $Q(s, a)$ used in the choice model.

1334 As described in original paper: “*The value of chosen features was updated based on the prediction*
 1335 *error scaled by a subject-specific update rate, while the value of unchosen features was decayed*
 1336 *toward 0 at a subject-specific decay rate* ([p461, Leong et al., 2017](#))”. We thereby updated this
 1337 feature values as:

$$1338 \quad \tilde{Q}(f, a_t) = \tilde{Q}(f, a_t) + \eta_{\text{attn}} \delta_t, \forall f \in \mathcal{F}(s_t) \quad (34)$$

1339 where $\delta_t = r_t - \tilde{Q}(s_t, a_t)$ and η_{attn} is the subject-specific update rate. This is very similar to Eq.
 1340 32, except for not modulated by the attention mechanism. For the unchosen features, we
 1341 decayed their value towards 0,

$$1342 \quad \tilde{Q}(f, a_t) = \tilde{Q}(f, a_t) + \epsilon \left(0 - \tilde{Q}(f, a_t) \right), \forall f \notin \mathcal{F}(s_t) \quad (35)$$

1343 where ϵ is the decay rate. Given that “...*the maximum feature value in each dimension was then*
 1344 *passed through a softmax function to obtain the predicted attention vector...* ([p461, Leong et al.,](#)
 1345 [2017](#))”, we defined the maximum feature value in each dimension, $d \in \{\text{shape, color,}$
 1346 $\text{appendage}\}$:

$$1347 \quad \xi(d) = \max_{f \in d} V(f) \quad (36)$$

1348 where $V(f) = \max_a Q(f, a)$, a popular method that has been widely used to approximate the
 1349 value of a feature marginalized over actions ([Sutton & Barto, 2018](#)). Thus, the attention can be
 1350 calculated as,

$$1351 \quad \phi(f) = \text{softmax}(\beta_{\text{attn}} \xi(d)), \quad \forall f \in d \quad (37)$$

1352 where β_{attn} is the inverse temperature.

1353 The five parameters $\{\eta, \beta, \eta_{\text{attn}}, \epsilon, \beta_{\text{attn}}\}$ of the ACL are: the learning rate for the value η ,
 1354 the inverse temperature β , the learning rate for the feature value in attention estimation η_{attn} ,
 1355 the decay rate of the feature value in attention estimation ϵ , the inverse temperature in attention
 1356 estimation β_{attn} .

1357

1358 4.4.10 L1PG, L2PG, and DCPG

1359 The three models have similar computational goals with the ECPG model, except that the
 1360 representation complexity $I^\psi(S; Z; \theta)$ terms in Eqs. 8 and 13 was respectively replaced by L1
 1361 norm ($\|\theta\|_1$), L2 norm ($\|\theta\|_2$), and decoder complexity $I^\rho(Z; A; \phi)$.

1362

1363 4.4.11 RNDPG: random regularizer policy gradient

1364 The RNDPG considers a commonly used regularizer in machine learning: injecting noise into the
 1365 encoder weights ([Noh et al., 2017](#)). Due to the injected noise $\varepsilon \in \mathbb{R}^{|F| \times |Z|}$ to the model, the
 1366 likelihood calculation should be turned to an expectation form,

$$1367 \quad \log \pi(a_t|s_t, \xi) = \log E_{\varepsilon \sim p(\varepsilon)}[\pi(a_t, \varepsilon|s_t, \xi)] \quad (38)$$

1368 Since our data lacks information regarding the type of noise the brain may have used in the
 1369 current experimental trial t , the most effective approach is to account for this uncertainty by
 1370 integrating over a presumed noise distribution. Accordingly, the update formula of the encoder
 1371 parameter should be updated as,

$$1372 \quad \max_{\theta} J(\theta) = \max_{\theta} E_{\varepsilon \sim p(\varepsilon)} \left[\sum_z \psi(z|s_t; \theta, \varepsilon) [\rho(a_t|z)(r(s_t, a_t) - b)] \right] \quad (39)$$

1373 Assume the noise follows a Gaussian distribution $N(0, \lambda^2)$, we can use the reparameterization
 1374 trick by first sampling ε from a standard Gaussian distribution $N(0, 1)$ and multiple by λ ,

$$1375 \quad \psi(z|s_t; \theta, \varepsilon) = \text{softmax}(\theta(f_t, z) + \lambda \varepsilon) \quad (40)$$

1376 Eq. 40 then has no closed-form solution, so we need to approximate it using the sampling method.

$$1377 \quad \max_{\theta} J(\theta) = \max_{\theta} \frac{1}{N_{\text{sample}}} \sum_{i=1}^{N_{\text{sample}}} \sum_z \psi(z|s_t; \theta, \varepsilon_i) [\rho(a_t|z)(r(s_t, a_t) - b)] \quad (41)$$

1378 Where $N_{\text{sample}} = 50$. Similarly, the update formula for decoder parameters is,

$$1379 \quad \max_{\phi} J(\phi) = \max_{\phi} \frac{1}{N_{\text{sample}}} \sum_{i=1}^{N_{\text{sample}}} \sum_z \rho(a_t|z; \phi) [\psi(z|s_t, \varepsilon_i)(r(s_t, a_t) - b)] \quad (42)$$

1380 Like other policy gradient model, the RNDPG model has three parameters $\{\alpha_\psi, \alpha_\rho, \lambda\}$:
 1381 the learning rate for the encoder α_ψ , the learning rate decoder α_ρ , and the standard deviation
 1382 of the injected noise λ .

1383

1384 4.5 Pretrain an encoder

1385 AE describes the phenomenon where generalization between two “superficially dissimilar”
 1386 stimuli increases after they have been paired with the same actions. To ensure their dissimilarity,
 1387 we selected stimuli that are easily distinguishable by human participants. We operationally

1388 defined this dissimilarity by setting a criterion: all four input stimuli must be classifiable with an
1389 accuracy of 99%. In order to accurately model human behaviors in the AE task, all models should
1390 undergo pretraining to achieve this level of discrimination accuracy.

1391 The key step of the pretraining is to search for an appropriate θ_0 . This is because all models'
1392 encoders were specially designed such that they can be initialized once the θ_0 is decided,

$$1393 \quad \bar{\psi}(z|s; \theta_0) = \begin{cases} \text{softmax}(\theta_0 \mathbf{I}(s, z)), & \text{for ECPG/CAPG} \\ \text{softmax}(\theta_0 \text{sim}(s, z)), & \text{for fECPG/fCAPG} \end{cases} \quad (43)$$

1394 The initial discrimination accuracy of the encoder can be quantified as follows:

$$1395 \quad \text{acc}(\theta_0) = \frac{1}{4} \sum_{i \in \{x, x', y, y'\}} \bar{\psi}(z = i | s = i; \theta_0) \quad (44)$$

1396 Given these constructs, we can search for an appropriate θ_0 by addressing the following
1397 objective,

$$1398 \quad \theta_0^* = \arg \min_{\theta_0} \frac{1}{2} (0.99 - \text{acc}(\theta_0))^2 \quad (45)$$

1399 Addressing this optimization objective, we initialized ECPG and CAPG model using $\theta_0^* =$
1400 5.232. To initialize the fECPG and fCAPG, we used $\theta_0^* = 1.459$ for the consistent and conflict case;
1401 $\theta_0^* = 1.329$ for the control case.

1402

1403 **4.6 Generate an action**

1404 A standard RL problem considers decision-making as sampling an action from the policy $\pi(a|s_t)$,
1405 a categorical distribution over the possible action space. In the AE problem, the possible action
1406 space varied from trial to trial. Participants were instructed to choose between a_1 and a_2 in one
1407 trial, while they were presented with a_3 and a_4 in another trial. To run both RL-base models in
1408 the AE task, we applied a technique called *invalid action masking* ([Huang & Ontañón, 2020](#)).

1409 The simplest masking is to add a large negative number ζ (in this thesis, $\zeta = -1e12$) to
1410 logits of the actions that are not presented in the current trial. That is, when an RLPG agent needs
1411 to choose between a_1 and a_2 , we can calculate its renormalized policy as,

$$1412 \quad \tilde{\pi}(a|s_t; \tilde{\phi}) = \text{softmax}(\tilde{\phi}(s_t, a)) \quad (46)$$

1413 where $\tilde{\phi}$ is,

1414
$$\tilde{\phi}(s, a) = \begin{cases} \phi(s, a), & \text{if } a \text{ is valid action at } t \\ \phi(s, a) + \zeta, & \text{if } a \text{ is invalid action at } t \end{cases} \quad (47)$$

1415 We then sampled from this renormalized policy to model human decision-making.

1416 For the ECPG and CAPG models, we masked and re-normalized the decoder for all
1417 representations z ,

1419
$$\tilde{\rho}(a|z; \tilde{\phi}) = \text{softmax}(\tilde{\phi}(z, a)) \quad (48)$$

1418

1420 **4.7 Perturbation-based feature importance**

1421 To investigate the importance of different features (color, shape, appendages) to humans, we
1422 adopted a *perturbation-based* measurement approach ([Greydanus et al., 2018](#); [Guo et al., 2021](#)),
1423 applied to our model fitted to human behavior. This method calculates the importance of each
1424 feature by theorizing that if an agent focuses heavily on a particular feature, then a minor
1425 perturbation in that feature might lead to significant changes in the output. This perturbation-
1426 based importance has been applied to extract measures of attention from large-scale deep
1427 reinforcement learning models in artificial intelligence, which was shown to be similar to human
1428 eye-tracking attention data ([Guo et al., 2021](#)).

1429 In the present work, we calculated the perturbation-based feature importance following
1430 the pseudo-algorithm 1 (**Extended Data 5.4**), modified based on the feature importance
1431 algorithm described in ([Fisher et al., 2019](#)).

1432

1433 **4.8 Model fitting**

1434 For each model, we estimated its free parameters separately for each subject, using all behavioral
1435 data from both the training and testing trials without cross-validation. This approach is consistent
1436 with many previous human learning studies, which are often structured with 2 to 4 parallel blocks
1437 due to various practical constraints (e.g. [Browning et al., 2015](#); [Daw et al., 2011](#); [Gershman, 2020](#);
1438 [Rac-Lubashevsky et al., 2023](#)). Given the insufficient number of blocks, these studies, including
1439 ours, do not meet the prerequisites for effective cross-validation.

1440 The parameters were estimated via maximum a posteriori (MAP):

$$1441 \quad \max_{\xi} \sum_{i=1}^N \log \pi(a_i|s_i, M, \xi) + \log p(\xi) \quad (49)$$

1442 where M refers to the model architecture, and ξ the model parameters. N is the number of trials
 1443 for each participant. s_i and a_i are the presented stimuli and human responses recorded on each
 1444 trial. We selected a very flat prior $p(\xi) = \text{Halfnorm}(0, 50)$ for all parameters with a range of
 1445 $(0, \infty)$ only to avoid extreme parameter values without biasing estimation. This prior is
 1446 uninformative yet ensures that parameter estimates remain within a reasonable range.
 1447 Parameters with a range of $(0, 1)$ used a uniform prior.

1448 Parameter estimation was performed using the BFGS algorithm, implemented with the
 1449 Python package `scipy.optimize.minimize`. For each participant, we ran the algorithm with 50 different
 1450 randomly chosen parameter initializations to avoid local minima in the non-convex landscape.

1451

1452 **4.9 Simulation**

1453 The parameter we are interested in is the simplicity degree (λ). We simulated the ECPG model's
 1454 learning and generalization behaviors by varying λ , while keeping the other two learning rate
 1455 parameters constant. In Experiment 1, the learning rate of the encoder was fixed at $\alpha_\psi = 40$,
 1456 and that of the decoder was fix at $\alpha_\rho = 4$. In Experiment 2, the learning rate of the encoder was
 1457 fixed at $\alpha_\psi = 8$, and that of the decoder was fix at $\alpha_\rho = 4$.

1458

1459 **4.10 Correlation between humans' and models' probe response**

1460 For each participants within a block, we calculated the frequency for each action as an estimation
 1461 of human probe policy. We applied the same method to the simulated data to obtain models'
 1462 probe policy. Subsequently, we computed the Spearman's correlation between human
 1463 participants and models based on the probability of selecting actions a_1 and a_3 . These two
 1464 actions sufficiently characterize a policy.

1465

1466 **5. Extended Data**

1467 **5.1 Model parameters**

1468 **Table S1: model parameters in Experiment 1 (mean \pm standard deviation)**

RLPG	$\alpha_\pi \in (0, \infty)$		
	0.567 \pm 0.630		
CAPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	
	4.376 \pm 6.519	1.553 \pm 1.753	
ECPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$
	14.983 \pm 14.881	2.336 \pm 2.926	0.112 \pm 0.277

1469 The reported values are from the 2.5% to 97.5% quantiles with extreme values removed.

1470

1471

Table S2: model parameters in Experiment 2 (mean \pm standard deviation)

fRLPG	$\alpha_\pi \in (0, \infty)$				
	0.567 \pm 0.630				
fCAPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$			
	4.376 \pm 6.519	1.553 \pm 1.753			
fECPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$		
	14.983 \pm 14.881	2.336 \pm 2.926	0.112 \pm 0.277		
LC	$\eta \in (0, 1)$	$\alpha \in (0, \infty)$	$\beta \in (0, \infty)$		
	0.108 \pm 0.232	1.637 \pm 5.515	127.218 \pm 416.861		
MA	$\alpha_Q \in (0, 1)$	$\alpha_{\text{assoc}} \in (0, 1)$	$\beta_{\text{assoc}} \in (0, \infty)$	$\beta \in (0, \infty)$	
	0.090 \pm 0.117	0.730 \pm 0.379	7.481 \pm 4.456	8.977 \pm 8.106	
ACL	$\eta \in (0, 1)$	$\beta \in (0, \infty)$	$\eta_{\text{attn}} \in (0, 1)$	$\epsilon \in (0, 1)$	$\beta_{\text{attn}} \in (0, \infty)$
	0.393 \pm 0.357	7.184 \pm 6.221	0.609 \pm 0.362	0.217 \pm 0.343	25.452 \pm 17.876
fL1PG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$		
	1.016 \pm 2.191	1.779 \pm 1.396	0.056 \pm 0.246		
fL2PG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$		
	1.479 \pm 2.279	1.579 \pm 1.185	0.107 \pm 0.671		
fRNDPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$		
	3.390 \pm 2.837	2.617 \pm 2.373	1.181 \pm 0.879		
fDCPG	$\alpha_\psi \in (0, \infty)$	$\alpha_\rho \in (0, \infty)$	$\lambda \in (0, \infty)$		

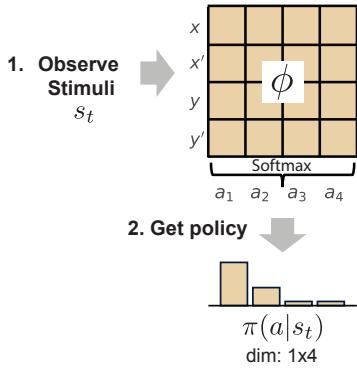
	30.620±43.514	0.027±0.095	0.012±0.038		
--	---------------	-------------	-------------	--	--

1472 The reported values are from the 2.5% to 97.5% quantiles with extreme values removed
1473

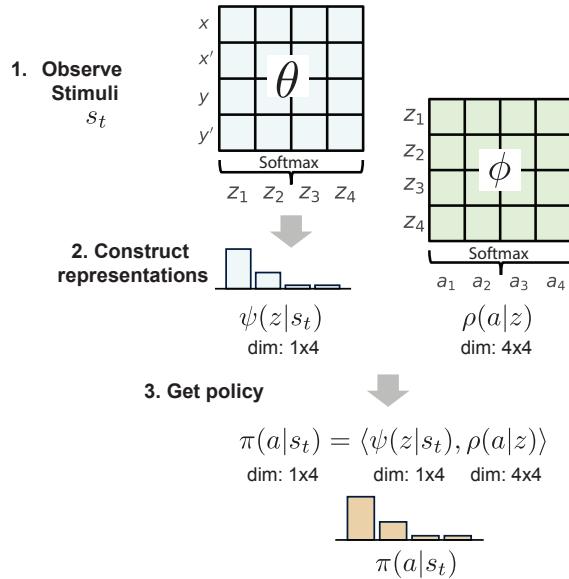
1474

5.2 A graphical illustration of the main models

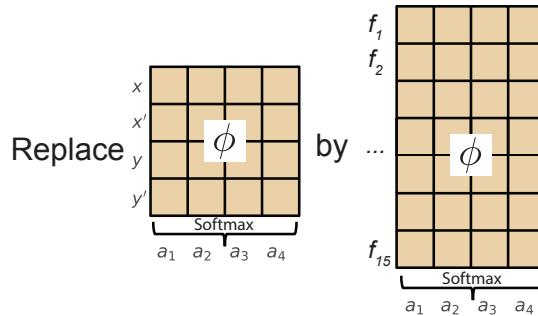
A The RLPG architecture



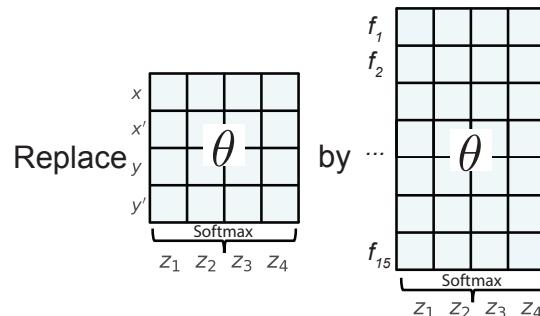
B The ECPG and CAPG architecture



C RLPG \rightarrow fRLPG



D ECPG \rightarrow fECPG



1475

1476

Figure S1 The semantic of the main model. $\langle \cdot, \cdot \rangle$ means dot product.

1477

1478

A. For each trial, the RLPG model observes a stimuli s_t and index the corresponding policy $\pi(a|s_t)$ for response.

1479

1480

B. The ECPG and the CAPG model needs to reconstruct the observed stimuli as an internal representation $\psi(z|s_t)$ and use the representation to derive a policy $\pi(a|s_t)$.

1481

1482

C. To build a feature-based RLPG, we only need to replace the 4(stimuli)-by-4(actions) policy with a 15 (features) by 4 (actions) policy.

1483

1484

D. To build a feature-based ECPG or CAPG, we only need to replace the 4(stimuli)-by-4(representations) encoder with a 15(features)-by-4(representations) encoder.

1485

1486

5.3 Representation evolution for all λ s

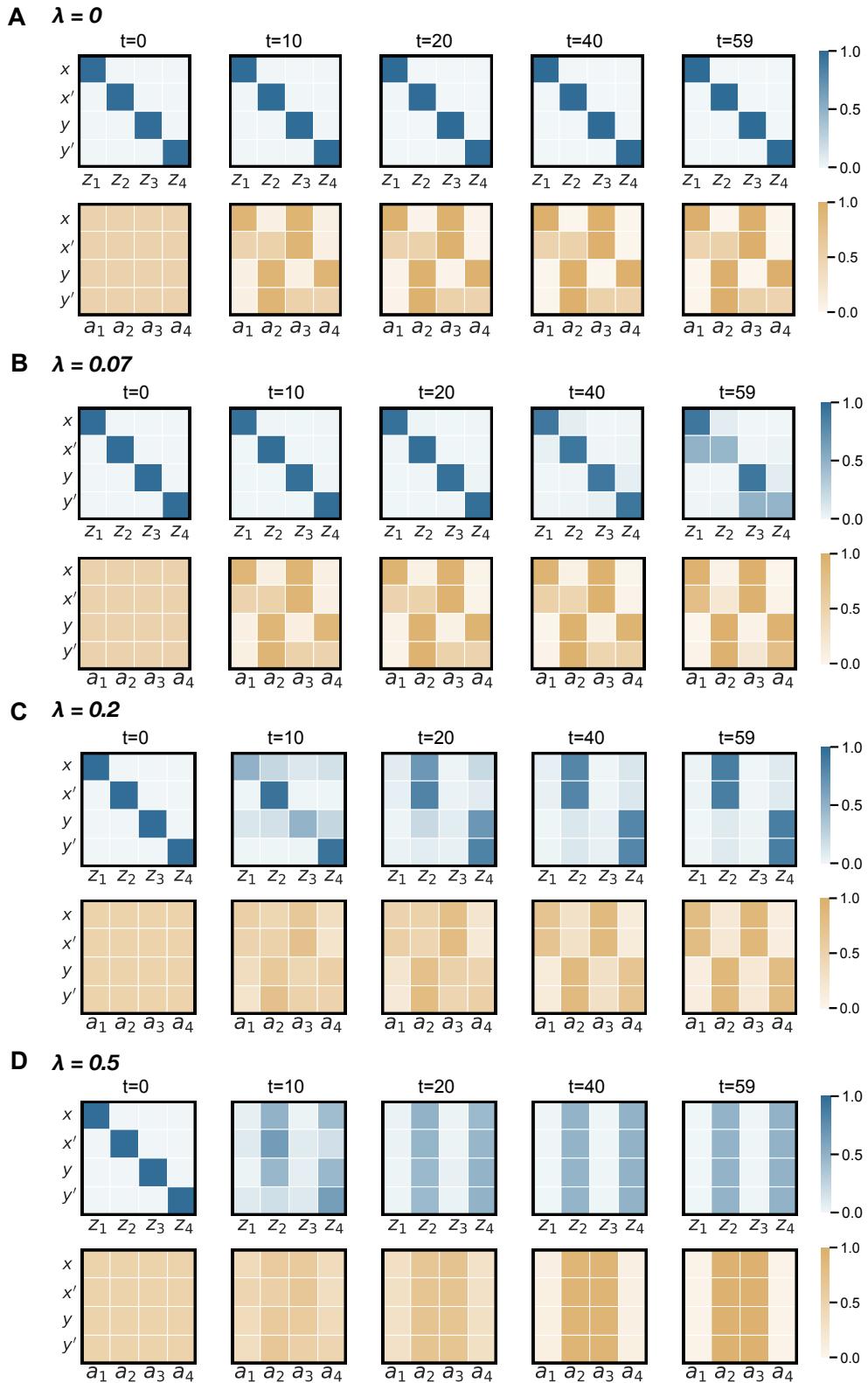


Figure S2 Encoders and Policies for ECPG model with different simplicity tendency, λ .

1490 **5.4 Pseudo algorithm for perturbation-based feature importance**

Algorithm 1 Calculate the perturbation-based feature importance of an encoder.

Require: Encoder ψ
Require: Feature embedding function \mathcal{F}
Require: Stimulus space $\{x, x', y, y'\}$
Require: Feature space {shape, color, appendage}
Require: Number of sample M

```

for  $d \in \{\text{shape, color, appendage}\}$  do
    for  $s \in \{x, x', y, y'\}$  do
         $f_{\text{orig}} = \mathcal{F}(s)$                                  $\triangleright$  Embed the stimulus
         $z_{\text{orig}} = \psi(f_{\text{orig}})$                    $\triangleright$  Encode the original stimulus
        for i=1:N do
             $f_{\text{pert}}^i = \text{change the one-hot code of feature } d \text{ to another value}$ 
             $z_{\text{pert}}^i = \mathcal{F}(f_{\text{pert}}^i)$                  $\triangleright$  Encode the perturbed stimulus
             $k^i = D_{\text{KL}}(z_{\text{orig}} \mid z_{\text{pert}})$            $\triangleright$  Compute the KL divergence
        end for
         $\text{imp}(s) = 1/N \sum_i^N k^i$        $\triangleright$  Importance, per stimulus per dimension
    end for
     $\text{IMPORTANT}(d) = 1/4 \sum_s \text{imp}(s)$        $\triangleright$  Importance, per dimension
end for
Return IMPORTANT

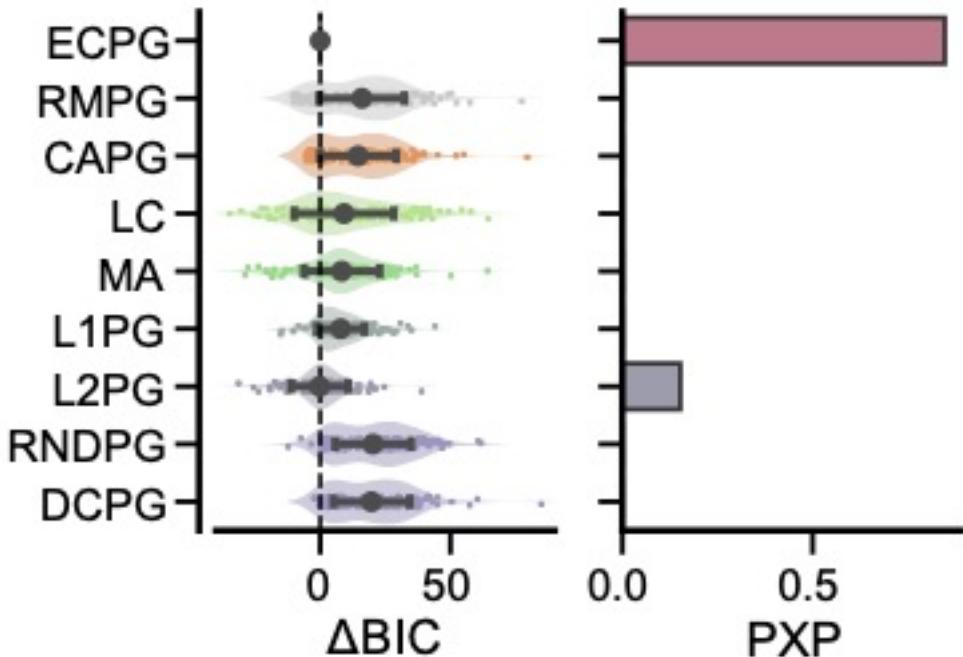
```

1491

1492

1493 **5.5 Model fitting in Experiment 1 for all models**

1494



1495

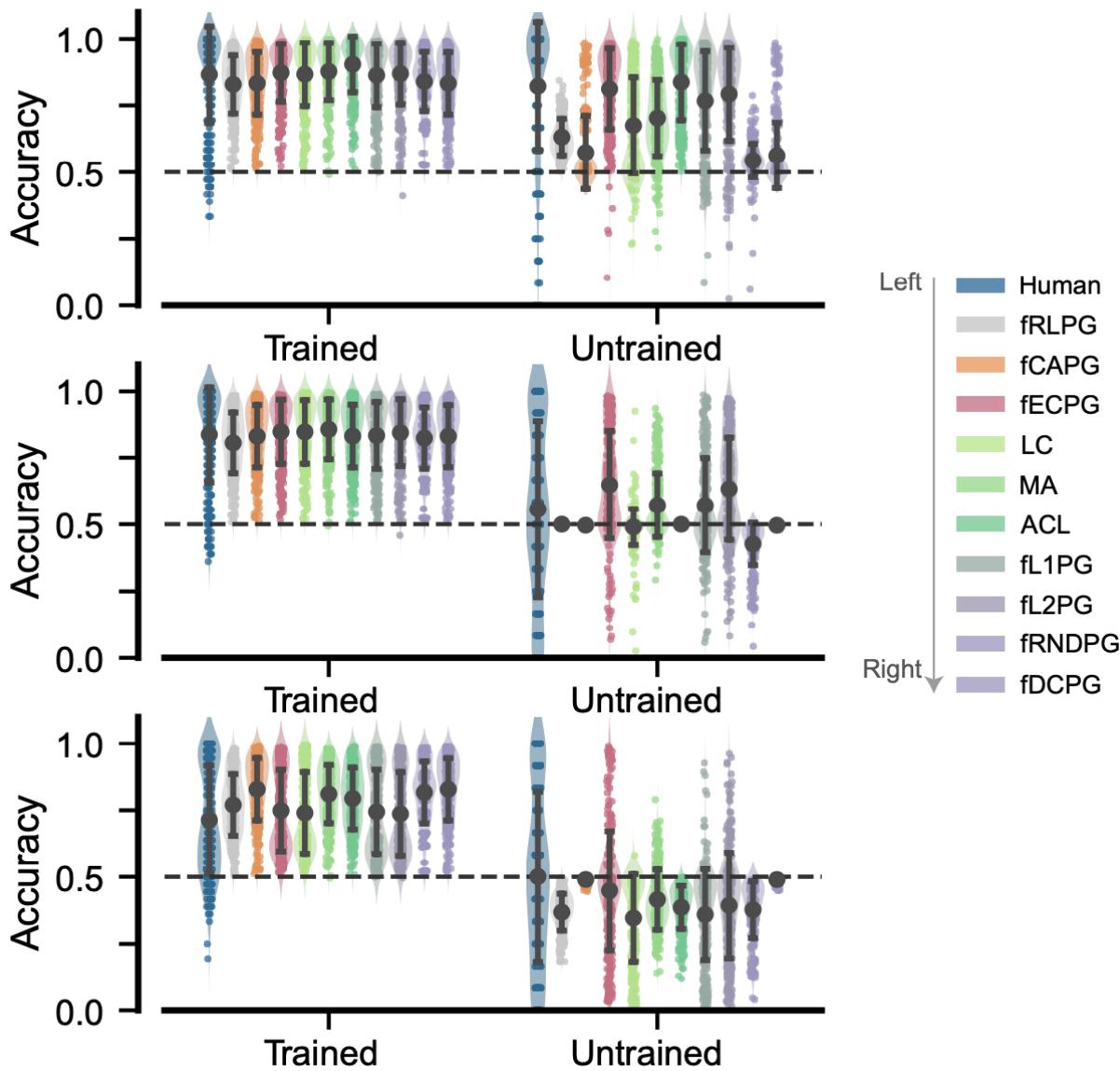
1496

1497 **Figure S3** Model fitting performance for all candidate models in Experiment 1. Error bars
1498 represent standard deviations. The ACL model, which must independently learn the value of
1499 each perceptual feature, was not applicable for Experiment 1. The L2PG model exhibits very
1500 close performance with our main ECPG model in terms of BIC. However, the PXP still
1501 significantly favors the ECPG model over the L2PG model (85% vs. 15%). This indicates that
1502 while the L2PG model is good at capturing certain behavioral patterns that are also discernible
1503 by other models, the ECPG model has a unique capability to identify specific behavioral
1504 patterns that other models fail to identify. See **Supplemental Note 2.3** for more details on the
1505 on the comparison between ECPG and L2PG. ECPG is the best model across the two datasets.
1506

1507

1508

5.6 Learning and generalization performance in the testing stage of Experiment 2



1509

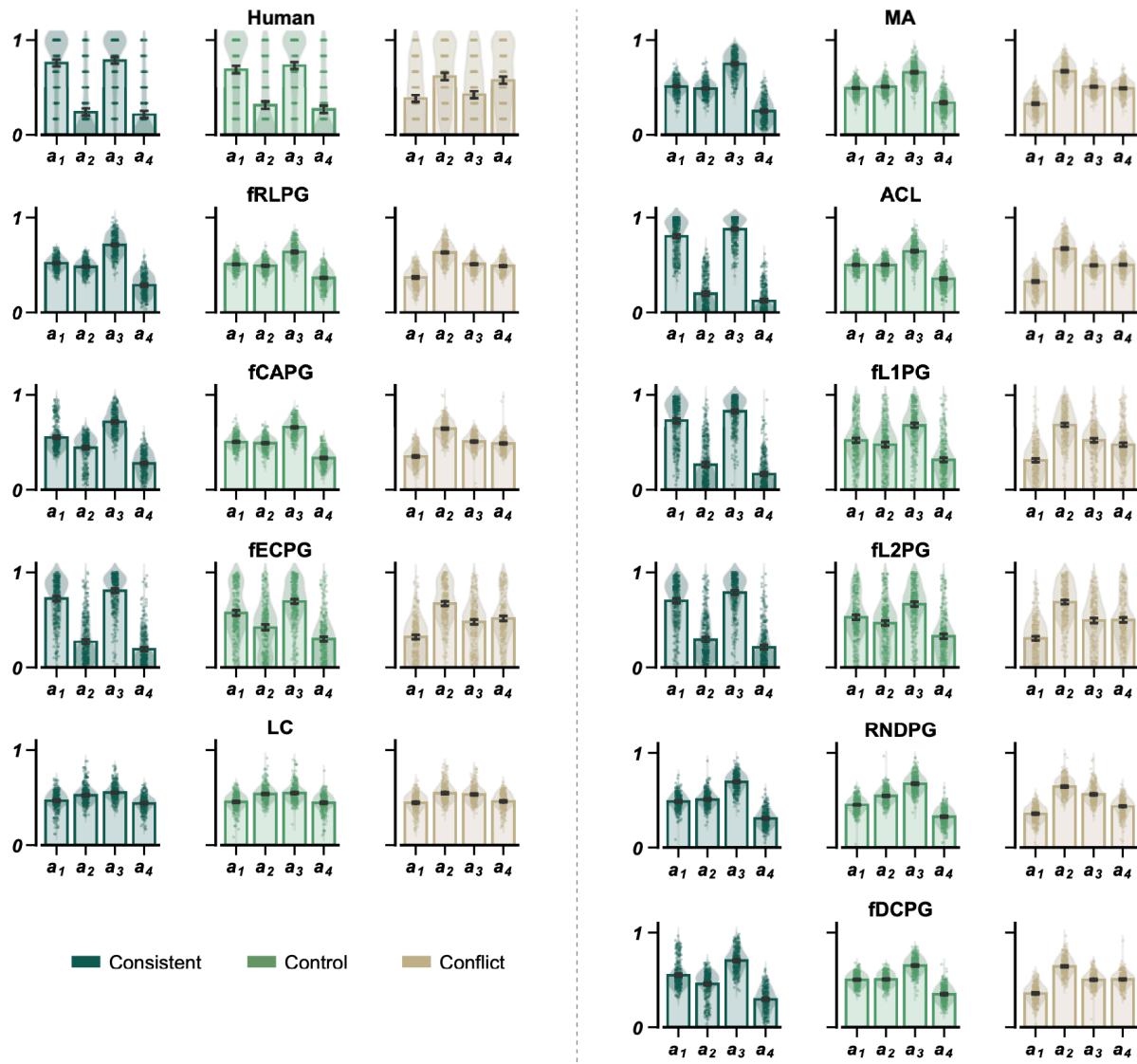
1510
1511
1512
1513
1514
1515
1516
1517
1518

Figure S4 Learning and generalization performance in the testing stage for all candidate models. Standard deviations are denoted by error bars. The models fRLPG, fCAPG, ACL, and fDCPG demonstrate a failure to generalize in the testing stage, suggesting a deficiency in their capacity for state abstraction. The two algorithmic-level models, LC and MA, which incorporate handcrafted state abstraction, show some degree of state abstraction capabilities but still fall short of human performance. In contrast, the fECPG model and its heuristic approximations, fL1PG and fL2PG, which aim to minimize representational complexity, effectively capture human learning and generalization capabilities.

1519

1520

5.7 Response to the probe stimuli



1521

1522

Figure S5: Responses to the probe stimuli for all candidate models. Error bars stand for a 95% confidence interval. The models fRLPG, fCAPG, LC, MA, and fDCPG cannot explain human probe responses, suggesting a deficiency in their capacity for rewarding feature extraction. The ACL model, equipped with a built-in feature extraction mechanism, only succeeds in consistent cases and fails in the other two. Finally, the fECPG model and its heuristic approximations, fL1PG and fL2PG, which aim to minimize representational complexity, successfully account for human response to the probe stimuli.

1523

1524

1525

1526

1527

1528

1529

1530

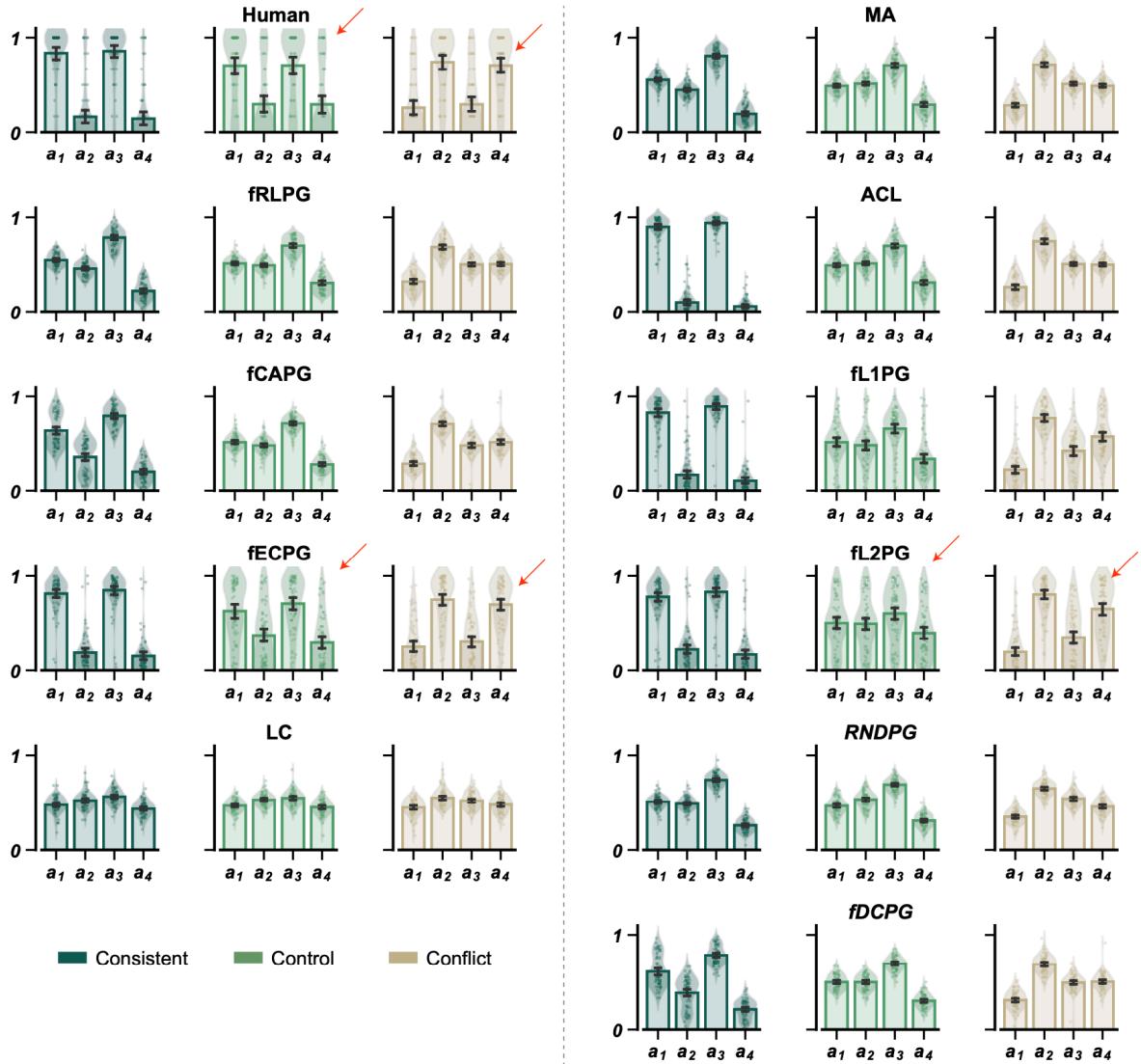
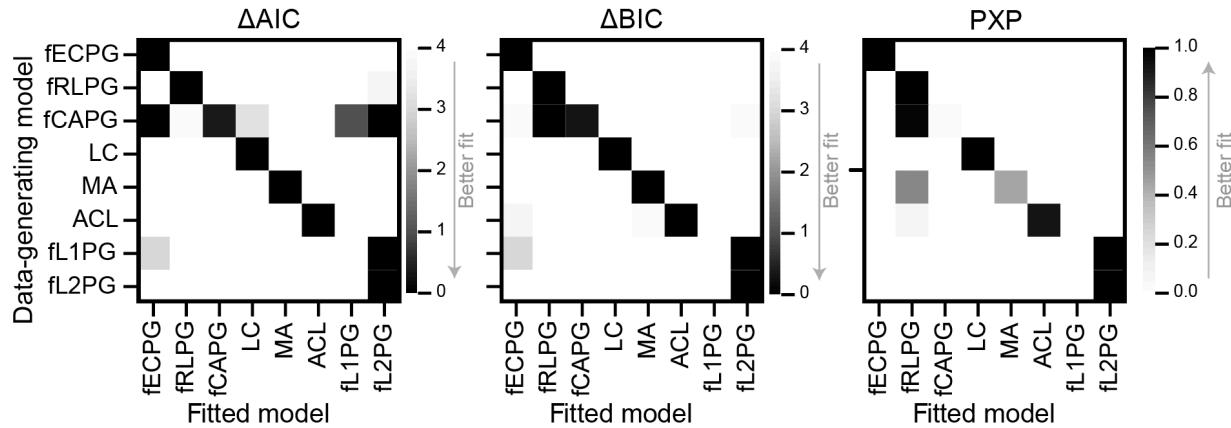


Figure S6: Responses to the probe stimuli for all candidate models for participants with better generalization. Out of 313 valid participants, we specifically selected 79 who has the top 25% for untrained accuracy to form a “high generalization” group. The fECPG model demonstrated a significantly higher fitting advantage over all other models, including the fL2PG. The red arrows highlight that the ECPG model’s alignment with human data is more precise in both control and conflict scenarios

1531
1532
1533
1534
1535
1536
1537
1538
1539

1540

5.8 The model recovery result



1541

1542 **Figure S6: Model recovery results for Experiment 2.** The AIC and BIC values are presented as
 1543 relative metrics, computed by subtracting the values from the best-fitting model (across
 1544 columns). The PXP metric is critical in the model recovery analysis. We sampled 40 participants
 1545 and used their fitted parameters to generate 10 simulated runs of behavioral data for
 1546 Experiment 2, with three blocks per run, each representing one experimental condition. This
 1547 resulted in 8 (models) \times 40 (parameter sets) \times 10 (samples) = 3,200 synthetic datasets. We
 1548 fitted all models to these synthetic datasets using the same fitting method as described
 1549 previously, except for the RNGPG and DCPG models due to their poor fitting performance. To
 1550 control for randomness, we averaged the fitting results (AICs and BICs) within each parameter
 1551 set over the 10 sample runs. The model recovery analysis reveals three key observations: first,
 1552 the ECPG model can be uniquely distinguished from the other models and does not falsely fit
 1553 well to simulated data from other models, demonstrating that its superior performance is due
 1554 to its accurate description of human behavior rather than greater expressiveness; second,
 1555 based on BIC and PXP, the CAPG model was always recovered as the RLPG model. This further
 1556 verifies that the success of the ECPG is not due to its cascade architecture but to its efficient
 1557 coding computational model; and third, the two norm-based regularizers (L1 and L2) cannot
 1558 be accurately differentiated from each other but can be separated from the ECPG model,
 1559 supporting our claim that the ECPG and the L2PG models capture different aspects of human
 1560 behavior. Overall, these findings indicate that our model recovery approach is sufficient to
 1561 confirm our conclusion that the ECPG model, with its augmented RL objective incorporating
 1562 efficient coding, best accounts for human learning and decision-making.
 1563