

1 **An active inference approach to modeling concept learning**

2 Ryan Smith,¹ Philipp Schwartenbeck,² Thomas Parr,² Karl J. Friston,²

3

4 ¹Laureate Institute for Brain Research, Tulsa, OK, USA

5 ²Wellcome Centre for Human Neuroimaging, Institute of Neurology, University

6 College London, WC1N 3BG, UK

7

8

9

10

11 Corresponding Author Information:

12

13 Ryan Smith

14 Laureate Institute for Brain Research

15 6655 S Yale Ave, Tulsa, OK 74136, USA

16 Email: rsmith@laureateinstitute.org

17

Abstract

18 Within computational neuroscience, the algorithmic and neural basis of concept
19 learning remains poorly understood. Concept learning requires both a type of
20 internal model expansion process (adding novel hidden states that explain new
21 observations), and a model reduction process (merging different states into one
22 underlying cause and thus reducing model complexity via meta-learning). Although
23 various algorithmic models of concept learning have been proposed within machine
24 learning and cognitive science, many are limited to various degrees by an inability to
25 generalize, the need for very large amounts of training data, and/or insufficiently
26 established biological plausibility. In this paper, we articulate a model of concept
27 learning based on active inference and its accompanying neural process theory, with
28 the idea that a generative model can be equipped with extra (hidden state or cause)
29 ‘slots’ that can be engaged when an agent learns about novel concepts. This can be
30 combined with a Bayesian model reduction process, in which any concept learning –
31 associated with these slots – can be reset in favor of a simpler model with higher
32 model evidence. We use simulations to illustrate this model’s ability to add new
33 concepts to its state space (with relatively few observations) and increase the
34 granularity of the concepts it currently possesses. We further show that it
35 accomplishes a simple form of ‘one-shot’ generalization to new stimuli. Although
36 deliberately simple, these results suggest that active inference may offer useful
37 resources in developing neurocomputational models of concept learning.

38 *Keywords:* Model Expansion; Structure Learning; Concepts; Computational
39 Neuroscience; Active Inference

41

Introduction

42 The ability to conceptualize and understand regularly observed patterns in
43 co-occurring feature observations is central to human cognition. For example, we do
44 not simply observe particular sets of colors, textures, shapes, and sizes – we also
45 observe *identifiable objects* such as, say, a ‘screwdriver’. If we were tool experts, we
46 might also recognize particular types of screwdrivers (e.g., flat vs. Phillip’s head),
47 designed for a particular use. This ability to recognize co-occurring features under
48 conceptual categories (as opposed to just perceiving sensory qualities; e.g., red,
49 round, etc.) is also highly adaptive. Only if we know an object is a screwdriver could
50 we efficiently infer that it affords putting certain structures together and taking
51 them apart; and only if we know the specific type of screwdriver could we efficiently
52 infer, say, the artefacts to use it on. Many concepts of this sort require experience-
53 dependent acquisition (i.e., learning).

54 From a computational perspective, the ability to acquire a new concept can
55 be seen as a type of Bayesian model comparison or structure learning (Botvinick,
56 Niv, & Barto, 2009; S. J. Gershman & Niv, 2010; MacKay & Peto, 1995; Salakhutdinov,
57 Tenenbaum, & Torralba, 2013; Tervo, Tenenbaum, & Gershman, 2016). Specifically,
58 concept acquisition can be cast as an agent learning (or inferring) that a new
59 hypothesis (referred to here as a hidden cause or state) should be added to the
60 internal or generative model with which she explains her environment, because
61 existing causes cannot account for new observations (e.g., an agent might start out
62 believing that the only tools are hammers and screwdrivers, but later learn that
63 there are also wrenches). In other words, the structure of the space of hidden causes

64 itself needs to expand to accommodate new patterns of observations. This model
65 expansion process is complementary to a process called Bayesian model reduction
66 (Karl Friston & Penny, 2011); in which the agent can infer that there is redundancy
67 in her model, and a model with fewer states or parameters provides a more
68 parsimonious (i.e. simpler) explanation of observations (KJ Friston, Lin, et al., 2017;
69 Schmidhuber, 2006). For example, in some instances it may be more appropriate to
70 differentiate between fish and birds as opposed to salmon, peacocks and pigeons.
71 This reflects a reduction in model complexity based on a particular feature space
72 underlying observations and thus resonates with other accounts of concept learning
73 as dimensionality reduction (Behrens et al., 2018; Stachenfeld, Botvinick, &
74 Gershman, 2016) – a topic we discuss further below.

75 A growing body of work in a number of domains has approached this
76 problem from different angles. In developmental psychology and cognitive science,
77 for example, probability theoretic (Bayesian) models have been proposed to account
78 for word learning in children and the remarkable human ability to generalize from
79 very few (or even one) examples in which both broader and narrower categorical
80 referents could be inferred (Kemp, Perfors, & Tenenbaum, 2007; Lake,
81 Salakhutdinov, & Tenenbaum, 2015; Perfors, Tenenbaum, Griffiths, & Xu, 2011; Xu &
82 Tenenbaum, 2007a, 2007b). In statistics, a number of nonparametric Bayesian
83 models, such as the “Chinese Room” process and the “Indian Buffet” process, have
84 been used to infer the need for model expansion (S. Gershman & Blei, 2012). There
85 are also related approaches in machine learning, as applied to things like Gaussian
86 mixture models (McNicholas, 2016).

87 These are the generative models that underwrite various clustering
88 algorithms that take sets of data points in a multidimensional space and divide them
89 into spatially separable clusters. While many of these approaches assume the
90 number of clusters is known in advance, various goodness-of-fit criteria may be
91 used to determine the optimal number. However, a number of approaches require
92 much larger amounts of data than humans do to learn new concepts (Geman,
93 Bienenstock, & Doursat, 1992; Hinton et al., 2012; LeCun, Bengio, & Hinton, 2015;
94 Lecun, Bottou, Bengio, & Haffner, 1998; Mnih et al., 2015). Many also require
95 corrective feedback to learn and yet fail to acquire sufficiently rich conceptual
96 structure to allow for generalization (Barsalou, 1983; Biederman, 1987; Feldman,
97 1997; Jern & Kemp, 2013; A. B. Markman & Makin, 1998; Osherson & Smith, 1981;
98 Ward, 1994; Williams & Lombrozo, 2010).

99 One potentially fruitful research avenue that has not yet been examined is to
100 explore concept learning from the perspective of Active Inference models based on
101 the free-energy principle (KJ Friston, 2010; KJ Friston et al., 2016; KJ Friston, Lin, et
102 al., 2017; KJ Friston, Parr, & de Vries, 2017). In this paper, we explore the potential
103 of this approach. In brief, we conclude that structure learning is an emergent
104 property of active inference (and learning) under generative models with ‘spare
105 capacity’; where spare or uncommitted capacity is used to expand the repertoire of
106 representations (Baker & Tenenbaum, 2014), while Bayesian model reduction (KJ
107 Friston, Lin, et al., 2017; Hobson & Friston, 2012) promotes generalization by
108 minimizing model complexity – and releasing representations to replenish ‘spare
109 capacity’.

110 In what follows, we first provide a brief overview of active inference. We then
111 introduce a model of concept learning (using basic and subordinate level animal
112 categories) to produce synthetic cognitive (semantic) processes associated with
113 adding new concepts to a state space and increasing the granularity of an existing
114 state space. We then establish the validity of this model using numerical analyses of
115 concept learning when repeatedly presenting a synthetic agent with different
116 animals characterized by different combinations of observable features. We will
117 demonstrate how particular approaches combining Bayesian model reduction and
118 expansion can reproduce successful concept learning without the need for
119 corrective feedback, and allow for generalization. We conclude with a brief
120 discussion of the implications of this work.

121

122 **An Active Inference model of concept learning**

123

124 **A primer on Active Inference**

125

126 Active Inference suggests that the brain is an inference machine that
127 approximates optimal probabilistic (Bayesian) belief updating across perceptual,
128 cognitive, and motor domains. Active Inference more specifically postulates that the
129 brain embodies an internal model of the world that is “generative” in the sense that
130 it can simulate the sensory data that it should receive if its model of the world is
131 correct. These simulated (predicted) sensory data can be compared to actual
132 observations, and differences between predicted and observed sensations can be

133 used to update the model. Over short timescales (e.g., a single observation) this
134 updating corresponds to inference (perception), whereas on longer timescales it
135 corresponds to learning (i.e., updating expectations about what will be observed
136 later). Another way of putting this is that perception optimizes beliefs about the
137 current state of the world, while learning optimizes beliefs about the relationships
138 between the variables that constitute the world. These processes can be seen as
139 ensuring the generative model (entailed by recognition processes in the brain)
140 remains an accurate model of the world that it seeks to regulate (Conant & Ashbey,
141 1970).

142 Active Inference casts decision-making in similar terms. Actions can be
143 chosen to resolve uncertainty about variables within a generative model (i.e.,
144 sampling from domains in which the model does not make precise predictions),
145 which can prevent anticipated deviations from predicted outcomes. In addition,
146 some expectations are treated as a fixed phenotype of an organism. For example, if
147 an organism did not continue to “expect” to observe certain amounts of food, water,
148 and shelter, then it would quickly cease to exist (McKay & Dennett, 2009) – as it
149 would not pursue those behaviors that fulfill these expectations (c.f. the ‘optimism
150 bias’ (Sharot, 2011)). Thus, a creature should continually seek out observations that
151 support – or are internally consistent with – its own continued existence. Decision-
152 making can therefore be cast as a process in which the brain infers the sets of
153 actions (policies) that would lead to observations consistent with its own survival-
154 related expectations (i.e., its “prior preferences”). Mathematically, this can be
155 described as selecting sequences of actions (policies) that maximize “Bayesian

156 model evidence” expected under a policy, where model evidence is the (marginal)
157 likelihood that particular sensory inputs would be observed under a given model.

158 In real-world settings, directly computing Bayesian model evidence is
159 generally intractable. Thus, some approximation is necessary. Active Inference
160 proposes that the brain computes a quantity called “variational free energy” that
161 provides a bound on model evidence, such that minimization of free energy
162 indirectly maximizes model evidence (this is exactly the same functional used in
163 machine learning where it is known as an evidence lower bound or ELBO). In this
164 case, decision-making will be approximately (Bayes) optimal if it infers (and enacts)
165 the policy that will minimize expected free energy (i.e., free energy with respect to a
166 policy, where one takes expected future observations into account). Technically,
167 expected free energy is the variational free energy averaged under the posterior
168 predictive density over policy-specific outcomes.

169 Expected free energy can be decomposed in different ways that reflect
170 uncertainty and prior preferences, respectively (e.g., epistemic and instrumental
171 affordance or ambiguity and risk). This formulation means that any agent that
172 minimizes expected free energy engages initially in exploratory behavior to
173 minimise uncertainty in a new environment. Once uncertainty is resolved, the agent
174 then exploits that environment to fulfil its prior preferences. The formal basis for
175 Active Inference has been thoroughly detailed elsewhere (KJ Friston, FitzGerald,
176 Rigoli, Schwartenbeck, & Pezzulo, 2017), and the reader is referred there for a full
177 mathematical treatment.

178 When the generative model is formulated as a partially observable Markov
179 decision process (a mathematical framework for modeling decision-making in cases
180 where some outcomes are under the control of the agent and others are not, and
181 where states of the world are not directly known but must be inferred from
182 observations), active inference takes a particular form. Here, the generative model is
183 specified by writing down plausible or allowable policies, hidden states of the world
184 (that must be inferred from observations), and observable outcomes, as well as a
185 number of matrices that define the probabilistic relationships between these
186 quantities (see right panel of figure 1).

Task Structure:

Learning animal concepts from observing co-occurring features

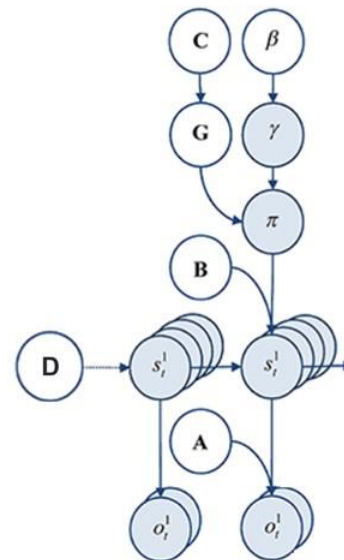


Time point 1
The agent observes the animal
(small, colourful, wings)



Time point 2
The agent reports the animal type

"PARAKEET"
"PARROT"
"PIGEON"
"BIRD"
...



187

188 Figure 1. Left: Illustration of the trial structure performed by the agent. At the first time
189 point, the agent is exposed to one of 8 possible animals that are each characterized by a
190 unique combination of visual features. At the 2nd time point, the agent would then report
191 which animal concept matched that feature combination. The agent could report a specific

192 category (e.g., pigeon, hawk, minnow, etc.) or a general category (i.e., bird or fish) if
193 insufficiently certain about the specific category. See the main text for more details. Right:
194 Illustration of the Markov decision process formulation of active inference used in the
195 simulations described in this paper. The generative model is heredeclared graphically, such
196 that arrows indicate dependencies between variables. Here observations (\mathbf{o}) depend on
197 hidden states (\mathbf{s}), as specified by the \mathbf{A} matrix, and those states depend on both previous
198 states (as specified by the \mathbf{B} matrix, or the initial states specified by the \mathbf{D} matrix) and the
199 policies ($\boldsymbol{\pi}$) selected by the agent. The probability of selecting a particular policy in turn
200 depends on the expected free energy (\mathbf{G}) of each policy with respect to the prior preferences
201 (\mathbf{C}) of the agent. The degree to which expected free energy influences policy selection is also
202 modulated by a prior policy precision parameter ($\boldsymbol{\gamma}$), which is in turn dependent on beta ($\boldsymbol{\beta}$)
203 –where higher values of beta promote more randomness in policy selection (i.e., less
204 influence of the differences in expected free energy across policies). For more details
205 regarding the associated mathematics, see (KJ Friston, Lin, et al., 2017; KJ Friston, Parr, et
206 al., 2017).

207

208 The ‘A’ matrix indicates which observations are generated by each
209 combination of hidden states (i.e., the likelihood mapping specifying the probability
210 that a particular set of observations would be observed given a particular set of
211 hidden states). The ‘B’ matrix is a transition matrix, indicating the probability that
212 one hidden state will evolve into another over time. The agent, based on the selected
213 policy, controls some of these transitions (e.g., those that pertain to the positions of
214 its limbs). The ‘D’ matrix encodes prior expectations about the initial hidden state
215 the agent will occupy. Finally, the ‘C’ matrix specifies prior preferences over
216 observations; it quantifies the degree to which different observed outcomes are
217 rewarding or punishing to the agent. In these models, observations and hidden
218 states can be factorized into multiple outcome *modalities* and hidden state *factors*.
219 This means that the likelihood mapping (the ‘A’ matrix) can also model the
220 interactions among different hidden states when generating outcomes
221 (observations). In what follows, we describe how this type of generative model was
222 specified to perform concept inference/learning.

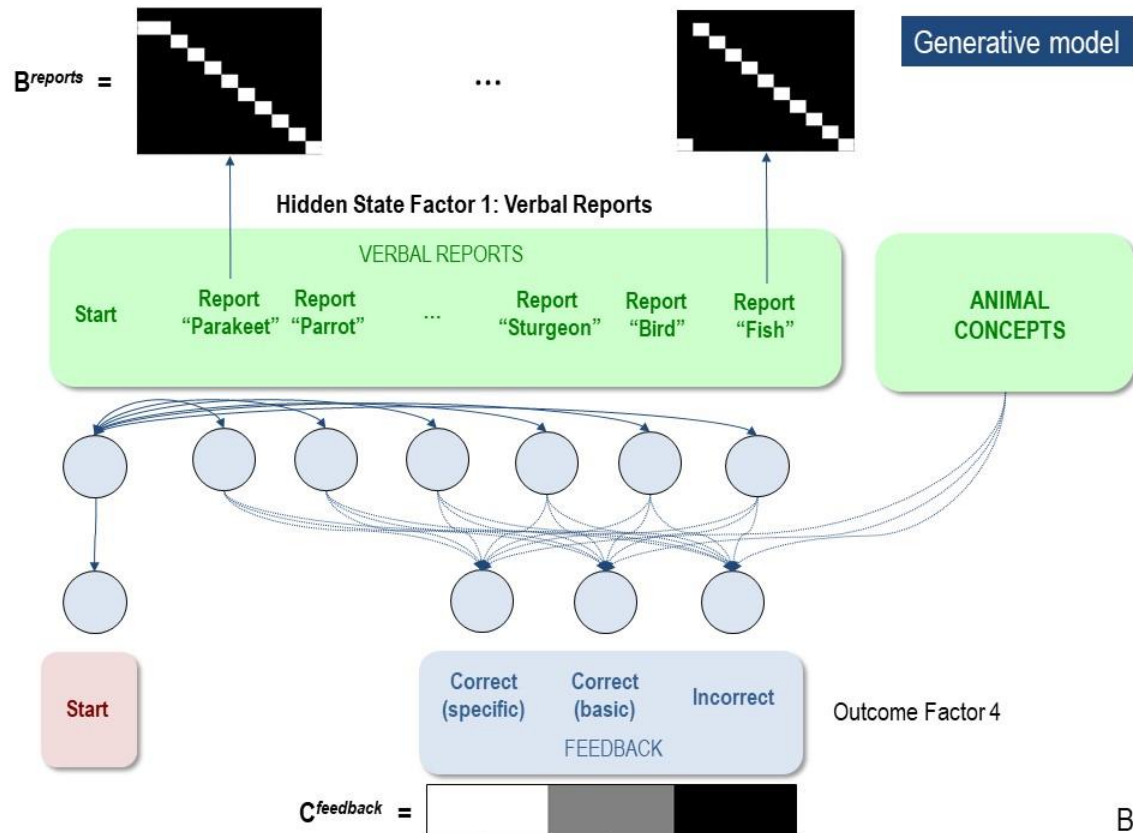
223

224 **A model of concept inference and learning**

225 To model concept inference, we constructed a simple task for an agent to
226 perform (see figure 1, left panel). In this task, the agent was presented with different
227 animals on different trials and asked to answer a question about the type of animal
228 that was seen. As described below, in some simulations the agent was asked to
229 report the type of animal that was learned previously; in other simulations, the
230 agent was instead asked a question that required conceptual generalization.

231 Crucially, to answer these questions the agent was required to observe different
232 animal features, where the identity of the animal depended on the combination of
233 features. There were three feature categories (size, color, and species-specific;
234 described further below) and two discrete time points in a trial (observe and
235 report).

236 To simulate concept learning (based on the task described above) we need to
237 specify an appropriate generative model. Once this model has been specified, one
238 can use standard (variational) message passing to simulate belief updating and
239 behavior in a biologically plausible way: for details, please see (KJ Friston,
240 FitzGerald, et al., 2017; KJ Friston, Parr, et al., 2017). In our (minimal) model, the
241 first hidden state factor included (up to) eight levels, specifying four possible types
242 of birds and four possible types of fish (Figure 2A). The outcome modalities
243 included: a feature space including two size features (big, small), two color features
244 (gray, colorful), and two species-differentiating features (wings, gills). The 'A' matrix
245 specified a likelihood mapping between features and animal concepts, such that



255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Figure 2. (A) Illustration of the first hidden state factor containing columns (levels) for 8 different animal concepts. Each of these 8 concepts generated a different pattern of visual feature observations associated with the outcome modalities of size, color, and species-specific features. The B matrix was an identity matrix, indicating that the animal being observed did not change within a trial (white = 1, black = 0). The A matrix illustrates the specific mapping from animal concepts to feature combinations. As depicted, each concept corresponded to a unique point in a 3-dimensional feature space. (B) illustration of the 2nd hidden state factor corresponding to the verbal reports the agent could choose in response to her observations. These generated feedback as to whether their verbal report was accurate with respect to a basic category report or a specific category report. As illustrated in the C matrix, the agent most preferred to be correct about specific categories, but least preferred being incorrect. Thus, reporting the basic categories was a safer choice if the agent was too uncertain about the specific identity of the animal.

270
271
272
273
274

The second hidden state factor was the agent's report. That this is assumed to factorise from the first hidden state factor means that there is no prior constraint that links the chosen report to the animal generating observations. The agent could report each of the eight possible specific animal categories, or opt for a less specific report of a bird or a fish. Only one category could be reported at any time. Thus, the

275 agent had to choose to report only bird vs. fish or to report a more specific category.
276 In other words, the agent could decide upon the appropriate level of coarse-graining
277 of her responses (figure 2B).

278 During learning trials, the policy space was restricted such that the agent
279 could not provide verbal reports or observe corrective feedback (i.e., all it could do
280 is “stay still” in its initial state and observe the feature patterns presented). This
281 allowed the agent to learn concepts in an unsupervised manner (i.e. without being
282 told what the true state was or whether it was correct or incorrect). After learning,
283 active reporting was enabled, and the ‘C’ matrix was set so that the agent preferred
284 to report correct beliefs. We defined the preferences of the agent such that she
285 preferred correctly reporting specific category knowledge and was averse to
286 incorrect reports. This ensured that she only reported the general category of bird
287 vs. fish, unless sufficiently certain about the more specific category.

288 In the simulations reported below, there were two time points in each trial of
289 categorisation or conceptual inference. At the first time point, the agent was
290 presented with the animals features, and always began in a state of having made no
291 report (the “start” state). The agent’s task was simply to observe the features, infer
292 the animal identity, and then report it (i.e., in reporting trials). Over 32 simulations
293 (i.e., 4 trials per animal), we confirmed that, if the agent already started out with full
294 knowledge of the animal concepts (i.e., a fully precise ‘A’ matrix), it would report the
295 specific category correctly 100% of the time. Over an additional 32 simulations, we
296 also confirmed that, if the agent was only equipped with knowledge of the
297 distinction between wings and gills (i.e., by replacing the rows in the ‘A’ matrix

298 corresponding to the mappings from animals to size and color with flat
299 distributions), it would report the generic category correctly 100% of the time but
300 would not report the specific categories.¹ This was an expected and straightforward
301 consequence of the generative model – but provides a useful example of how agents
302 trade off preferences and different types of uncertainty.

303

304 **Simulating concept learning and the acquisition of expertise**

305

306 Having confirmed that our model could successfully recognize animals if
307 equipped with the relevant concepts (i.e., likelihood mappings) – we turn now to
308 concept learning.

309

310 **Concept acquisition**

311 We first examined our model’s ability to acquire concept knowledge in two
312 distinct ways. This included 1) the agent’s ability to “expand” (i.e., fill in an unused
313 column within) its state space and add new concepts, and 2) the agent’s ability to
314 increase the granularity of its conceptual state space and learn more specific
315 concepts, when it already possessed broader concepts.

316

317 *Adding Concepts*

¹ However, "risky" reporting behavior could be elicited by manipulating the strengths of the agent's preferences such that she placed a very high value on reporting specific categories correctly (i.e., relative to how much it disliked reporting incorrectly).

318 To assess whether our agent could expand her state space by acquiring a new
319 concept, we first set one column of the previously described model's 'A' matrix
320 (mapping an animal concept to its associated features) to be a uniform distribution²;
321 creating an imprecise likelihood mapping for one concept – essentially, that concept
322 predicted all features with nearly equal probability. Here, we chose sturgeon (large,
323 gray, gills) as the concept for which the agent had no initial knowledge (see Figure
324 3A, right-most column of left-most 'pre-learning' matrix). We then generated 2000
325 observations based on the outcome statistics of a model with full knowledge of all
326 eight animals (the "generative process"), to test whether the model could learn the
327 correct likelihood mapping for sturgeon (note: this excessive number of
328 observations was used for consistency with later simulations, in which more
329 concepts have to be learned, and also to evaluate how performance improved as a
330 function of the number of observations the agent was exposed to; see figure 3B).

331 In these simulations, learning was implemented via updating (concentration)
332 parameters for the model's 'A' matrix after each trial. For details of these free energy
333 minimizing learning processes, please see (KJ Friston et al., 2016) as well as the left
334 panel of Figure 8 and the associated legend further below. An intuitive way to think
335 about this belief updating process is that the strength of association between a
336 concept and an observation is quantified simply by counting how often they are
337 inferred to co-occur. This is exactly the same principle that underwrites Hebbian
338 plasticity and long-term potentiation (Brown, Zhao, & Leung, 2010). Crucially,

² To break the symmetry of the uniform distribution, we added small amounts of Gaussian noise (with a variance of .001) to avoid getting stuck in local free energy minima during learning.

339 policies were restricted during learning, such that the agent could not select
340 reporting actions; thus, learning was driven entirely by repeated exposure to
341 different feature combinations. We evaluated successful learning in two ways. First,
342 we compared the ‘A’ matrix learned by the model to that of the generative process.
343 Second, we disabled learning after various trial numbers (i.e., such that
344 concentration parameters no longer accumulated) and enabled reporting. We then
345 evaluated reporting accuracy with 20 trials for each of the 8 concepts.

346 As shown in Figure 3A, the ‘A’ matrix (likelihood) mapping – learned by the
347 agent – and the column for sturgeon in particular, strongly resembled that of the
348 generative process. When first evaluating reporting, the model was 100 % accurate
349 across 20 reporting trials, when exposed to a sturgeon (reporting accuracy when
350 exposed to each of the other animals also remain at 100%) and first reached this
351 level of accuracy after around 50 exposures to all 8 animals (with equal probability)
352 (figure 3B). The agent also always chose to report specific categories (i.e., it never
353 chose to only report bird or fish). Model performance was stable over 8 repeated
354 simulations.

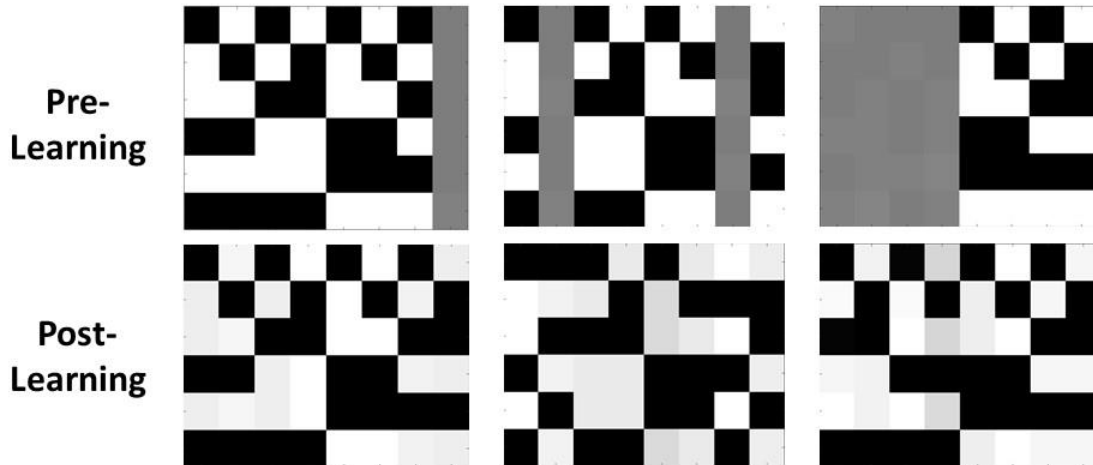
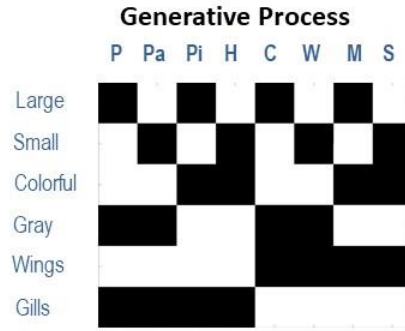
355 Crucially, during learning, the agent was not told which state was generating
356 its observations. This meant that it had to solve both an inference and a learning
357 problem. First, it had to infer whether a given feature combination was better
358 explained by an existing concept, or by a concept that predicts features uniformly. In
359 other words, it had to decide that the features were sufficiently different – from
360 things it had seen before – to assign it a new hypothetical concept. Given that a novel
361 state is only inferred when another state is not a better explanation, this precludes

362 learning 'duplicate' states that generate the same patterns of observations. The
363 second problem is simpler. Having inferred that these outcomes are caused by
364 something new, the problem becomes one of learning a simple state-outcome
365 mapping through accumulation of Dirichlet parameters.

366 To examine whether this result generalized, we repeated these simulations
367 under conditions in which the agent had to learn more than one concept. When the
368 model needed to learn one bird (parakeet) and one fish (minnow), the model was
369 also able to learn the appropriate likelihood mapping for these 2 concepts (although
370 note that, because the agent did not receive feedback about the state it was in during
371 learning, the new feature mappings were often not assigned to the same columns as
372 in the generative process; see figure 3A). Reporting also reached 100% accuracy,
373 but required a notably greater number of trials. Across 8 repeated simulations, the
374 mean accuracy reached by the model after 2000 trials was 98.75% (SD = 2%).

A

Adding Concepts: Categories as Columns

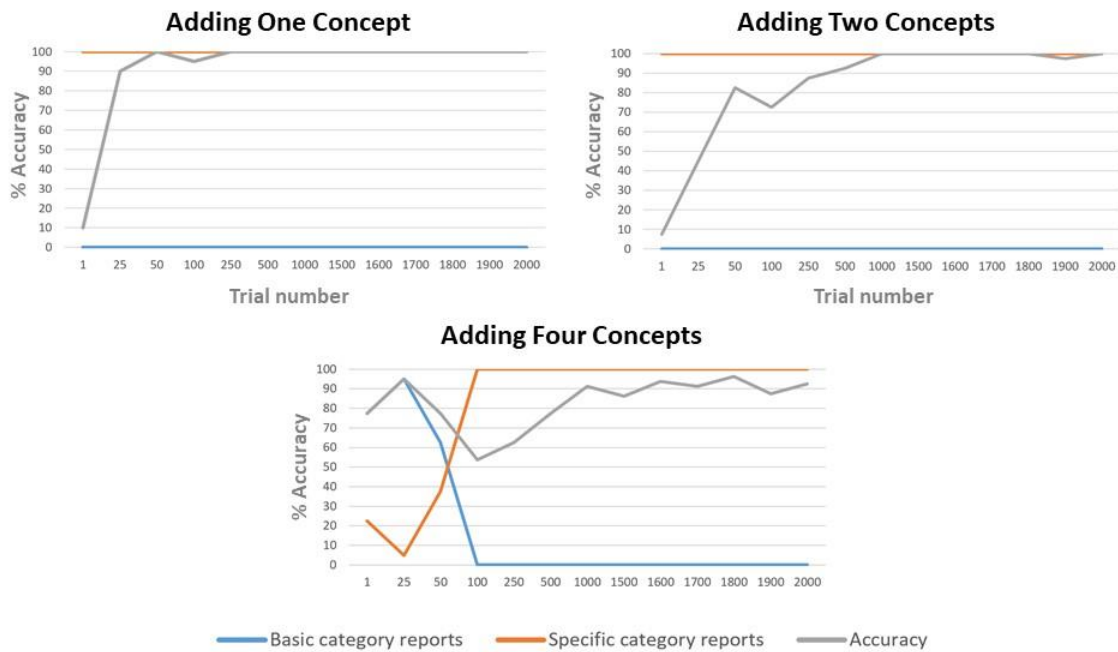


375

B

Adding Concepts:

Reporting Accuracy as a Function of Trial Number (without feedback)



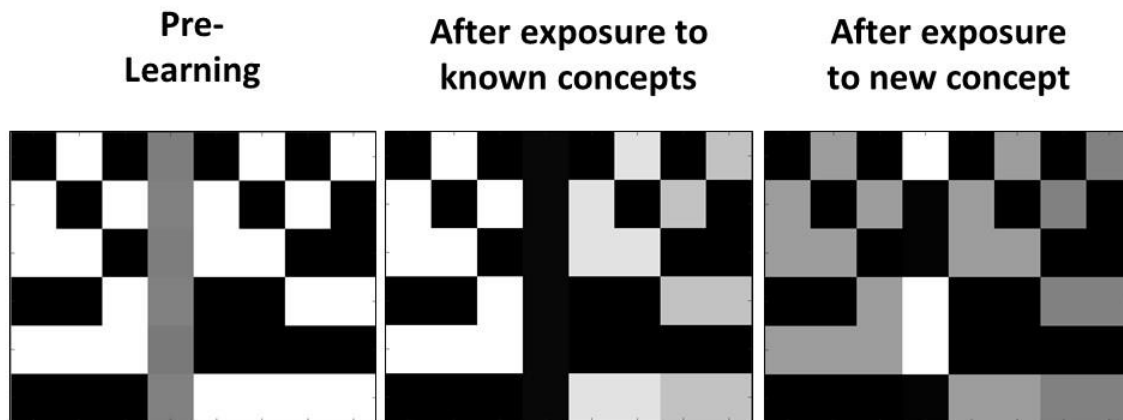
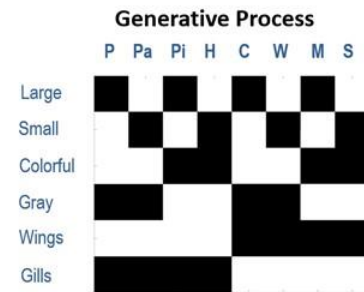
376

377 Figure 3. (A) illustration of representative simulation results in which the agent successfully
378 learned 1, 2, or 4 new animal concept categories with no prior knowledge beforehand. The
379 generative process is shown in the upper right, illustrating the feature combinations to be
380 learned. Pre-learning, either 1, 2 or 4 columns in the likelihood mapping began as a flat
381 distribution with a slight amount of Gaussian noise. The agent was and then provided with
382 2000 observations of the 8 animals with equal probability. Crucially, the agent was
383 prevented from providing verbal reports during these 2000 trials and thus did not receive
384 feedback about the true identity of the animal. Thus learning was driven completely by
385 repeated exposure in an unsupervised manner. Also note that, while the agent was
386 successful at learning the new concepts, it did not always assign the new feature patterns to
387 the same columns as illustrated in the generative process. This is to be expected given that
388 the agent received no feedback about the true hidden state that generated her observations.
389 (B) illustration of how reporting accuracy, and the proportion of basic category and specific
390 category responses, changed as a function of repeated exposures. This was accomplished by
391 taking the generative model at a number of representative trials and then testing it with 20
392 observations of each animal in which reporting was enabled. As can be seen, maximal
393 accuracy was achieved much more quickly when the agent had to learn fewer concepts.
394 When it had learned 4 concepts, it also began by reporting the general categories and then
395 subsequently became sufficiently confident to report the more specific categories.
396

397 When the model needed to learn all 4 birds, performance varied somewhat
398 more when the simulations were repeated. The learned likelihood mappings after
399 2000 trials always resembled that of the generative process, but with variable levels
400 of precision; notably, the model again assigned different concepts to different
401 columns relative to the generative process, as would be expected when the agent is
402 not given feedback about the state she is in. Over 8 repeated simulations, the model
403 performed well in 6 (92.50 % – 98.8 % accuracy) and failed to learn one concept in
404 the other 2 (72.50 % accuracy in each) due to overgeneralization (e.g., mistaking
405 parrot for Hawk in a majority of trials; i.e., the model simply learned that there are
406 large birds). Figure 3A and 3B illustrate representative results when the model was
407 successful (note: the agent never chose to report basic categories in the simulations
408 where only 1 or 2 concepts needed to be learned).

409 To further assess concept learning, we also tested the agent's ability to
410 successfully avoid state duplication. That is, we wished to confirm that the model
411 would only learn a new concept if actually presented with a new animal for which it
412 did not already have a concept. To do so, we equipped the model with knowledge of
413 seven out of the eight concept categories, and then repeatedly exposed it only to the
414 animals it already knew over 80 trials. We subsequently exposed it to the eighth
415 animal (Hawk) for which it did not already have knowledge over 20 additional
416 trials. As can be seen in figure 4, the unused concept column was not engaged during
417 the first 80 trials (bottom left and middle). However, in the final 20 trials, the agent
418 correctly inferred that her current conceptual repertoire was unable to explain her
419 new pattern of observations, leading the unused concept column to be adumbrated
420 and the appropriate state-observation mapping to be learned (bottom right). We
421 repeated these simulations under conditions in which the agent already had
422 knowledge of six, five, or four concepts. In all cases, we observed that unused
423 concept columns were never engaged inappropriately.
424

Avoiding State Duplication



425
426 Figure 4. Illustration of representative simulation results when the agent had to avoid
427 inappropriately learning a new concept (i.e., avoid state duplication) after only being
428 exposed to animals for which it already had knowledge. Here the agent began with prior
429 knowledge about seven concept categories, and was also equipped with an eighth column
430 that could be engaged to learn a new concept category (bottom left). The agent was then
431 presented with several instances of each of the seven animals that she already knew (80
432 trials in total). In this simulation, the agent was successful in assigning each stimulus to an
433 animal concept she had already acquired, and did not engage the unused concept 'slot'
434 (bottom middle). Finally, the agent was presented with a new animal (a hawk) that she did
435 not already know over 20 trials. In this case, the agent successfully engaged the additional
436 column (i.e., she inferred that none of the concepts she possessed could account for her new
437 observations), and learn the correct state-observation mapping (bottom right).
438

439 Crucially, these simulations suggest that adaptive concept learning needs to
440 be informed by existing knowledge about other concepts, such that a novel concept
441 should only be learned if observations cannot be explained with existing conceptual
442 knowledge. Here, this is achieved via the interplay of inference and learning, such
443 that agents initially have to infer whether to assign an observation to an existing

444 concept, and only if this is not possible an ‘open slot’ is employed to learn about a
445 novel concept.

446

447 *Increasing granularity*

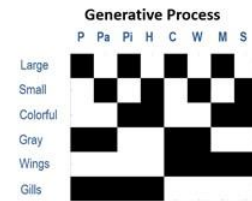
448 Next, to explore the model’s ability to increase the granularity of its concept
449 space, we first equipped the model with only the distinction between birds and fish
450 (i.e., the rows of the likelihood mapping corresponding to color and size features
451 were flattened in the same manner described above). We then performed the same
452 procedure used in our previous simulations. As can be seen in Figure 5A (bottom
453 left), the ‘A’ matrix learned by the model now more strongly resembled that of the
454 generative process. Figure 5A (bottom) also illustrates reporting accuracy and the
455 proportion of basic and specific category reports as a function of trial number. As
456 can be seen, the agent initially only reported general categories, and became
457 sufficiently confident to report specific categories after roughly 50 – 100 trials, but
458 its accuracy increased gradually over the next 1000 trials (i.e., the agent reported
459 specific categories considerably before its accuracy improved). Across 8 repeated
460 simulations, the final accuracy level reached was between 93% – 98% in 7
461 simulations, but she failed to learn one concept in the 8th case, with 84.4% overall
462 accuracy (i.e., a failure to distinguish between pigeon and parakeet, and therefore
463 only learned a broader category of “small birds”).

464 To assess whether learning basic categories first was helpful in subsequently
465 learning specific categories, we also repeated this simulation without any initial
466 knowledge of the basic categories. As exemplified in figure 5A and 5B, the model

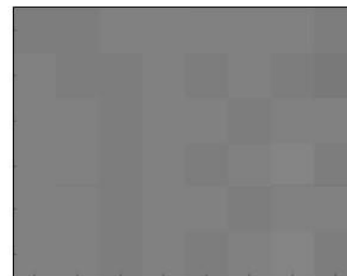
467 tended to perform reasonably well, but most often learned a less precise likelihood
468 mapping and reached a lower reporting accuracy percentage after 2000 learning
469 trials (across 8 repeated simulations: mean = 81.21%, SD = 6.39%, range from
470 68.80% – 91.30%). Thus, learning basic concept categories first appeared to
471 facilitate learning more specific concepts later.

A

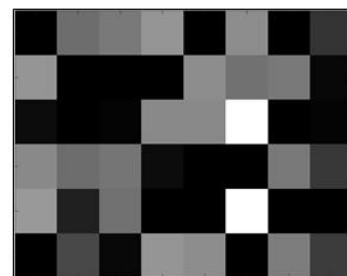
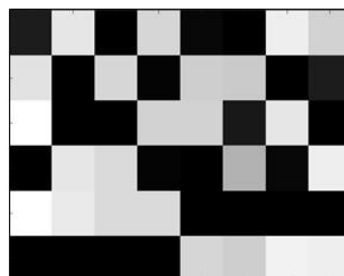
Increasing Granularity: Features as Rows



Pre-
Learning



Post-
Learning

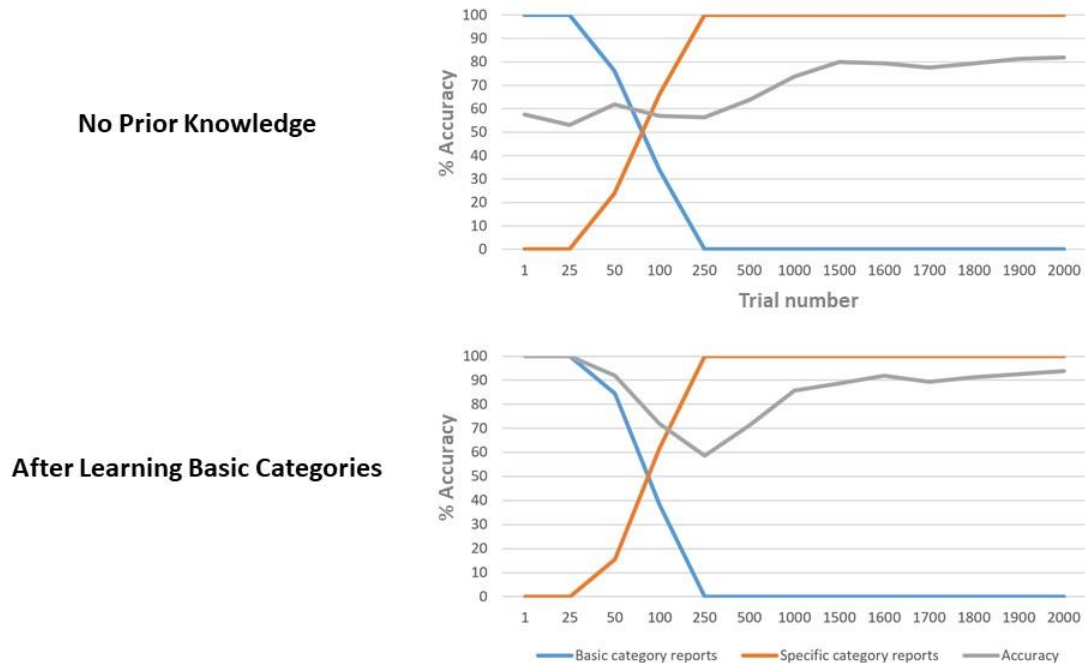


472
473

B

Increasing Granularity:

Reporting Accuracy as a Function of Trial Number (without feedback)



474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

Figure 5. (A, left) Illustration of representative simulation results when the agent had to learn to increase the granularity of its concept space. Here the agent began with prior knowledge about the basic concept categories (i.e., it had learned the broad categories of “bird” and “fish”) but had not learned the feature patterns (i.e., rows) that differentiate different types of birds and fish. Post learning (i.e., after 2000 exposures), the agent did successfully learn all of the more granular concept categories, although again note that specific concepts were assigned to different columns than depicted in the generative process due to the unsupervised nature of the learning. (A, right) illustration of the analogous learning result when the agent had to learn all 8 specific categories without prior knowledge of the general categories. Although moderately successful, learning tended to be more difficult in this case. (B) Representative plots of reporting accuracy in each of the 2 learning conditions as a function of the number of exposures. As can be seen, when the model starts out with prior knowledge about basic categories, it slowly become sufficiently confident to start reporting the more specific categories, and its final accuracy is high. In contrast, while the agent that did not start out with any prior knowledge of the general categories also grew confident in reporting specific categories over time, its final accuracy levels tended to be lower. In both cases, the agent began reporting specific categories before it achieved significant accuracy levels, therefore showing some initial overconfidence.

495 Overall, these findings provide a proof of principle that this sort of active
496 inference scheme can add concepts to its state space in an unsupervised manner
497 (i.e., without feedback) based purely on (expected) free energy minimization. In this
498 case, it was able to accomplish this starting from a completely uninformative
499 likelihood distribution. It could also learn more granular concepts after already
500 acquiring more basic concepts, and our results suggest that learning granular
501 concepts may be facilitated by first learning basic concepts (e.g., as in currently
502 common educational practices).

503 The novel feature of this generative model involved ‘building in’ a number of
504 “reserve” hidden state levels. These initially had uninformative likelihood mappings;
505 yet, if a new pattern of features was repeatedly observed, and the model could not
506 account for this pattern with existing (informative) state-observation mappings,
507 these additional hidden state levels could be engaged to improve the model’s
508 explanatory power. This approach therefore accommodates a simple form of model
509 expansion.

510

511 **Integrating model expansion and reduction**

512

513 We next investigated ways in which model expansion could be combined with
514 Bayesian model reduction (KJ Friston, Lin, et al., 2017) – allowing the agent to adjust
515 her model to accommodate new patterns of observations, while also precluding
516 unnecessary conceptual complexity (i.e., over-fitting). To do so, we again allowed
517 the agent to learn from 2000 exposures to different animals as described in the

518 previous section – but also allowed the model to learn its ‘D’ matrix (i.e., accumulate
519 concentration parameters reflecting prior expectations over initial states). This
520 allowed an assessment of the agent’s probabilistic beliefs about which hidden state
521 factor levels (animals) it had been exposed to. In different simulations, the agent
522 was only exposed to some animals and not others. We then examined whether a
523 subsequent model reduction step could recover the animal concepts presented
524 during the simulation; eliminating those concepts that were unnecessary to explain
525 the data at hand. The success of this 2-step procedure could then license the agent to
526 “reset” the unnecessary hidden state columns after concept acquisition, which
527 would have accrued unnecessary likelihood updates during learning. Doing so
528 would allow the optimal ability for those “reserve” states to be appropriately
529 engaged, if and when the agent was exposed to truly novel stimuli.

530 The 2nd step of this procedure was accomplished by applying Bayesian model
531 reduction to the ‘D’ matrix concentration parameters after learning. This is a form of
532 post-hoc model optimization (K. J. Friston et al., 2016; Karl Friston, Parr, & Zeidman,
533 2018) that rests upon estimation of a ‘full’ model, followed by analytic computation
534 of the evidence that would have been afforded to alternative models (with
535 alternative, ‘reduced’, priors) had they been used instead. Mathematically, this
536 procedure is a generalization of things like automatic relevance determination (Karl
537 Friston, Mattout, Trujillo-Barreto, Ashburner, & Penny, 2007; Wipf & Rao, 2007) or
538 the use of the Savage Dickie ratio in model comparison (Cornish & Littenberg,
539 2007). It is based upon straightforward probability theory and, importantly, has a
540 simple physiological interpretation; namely, synaptic decay and the elimination of

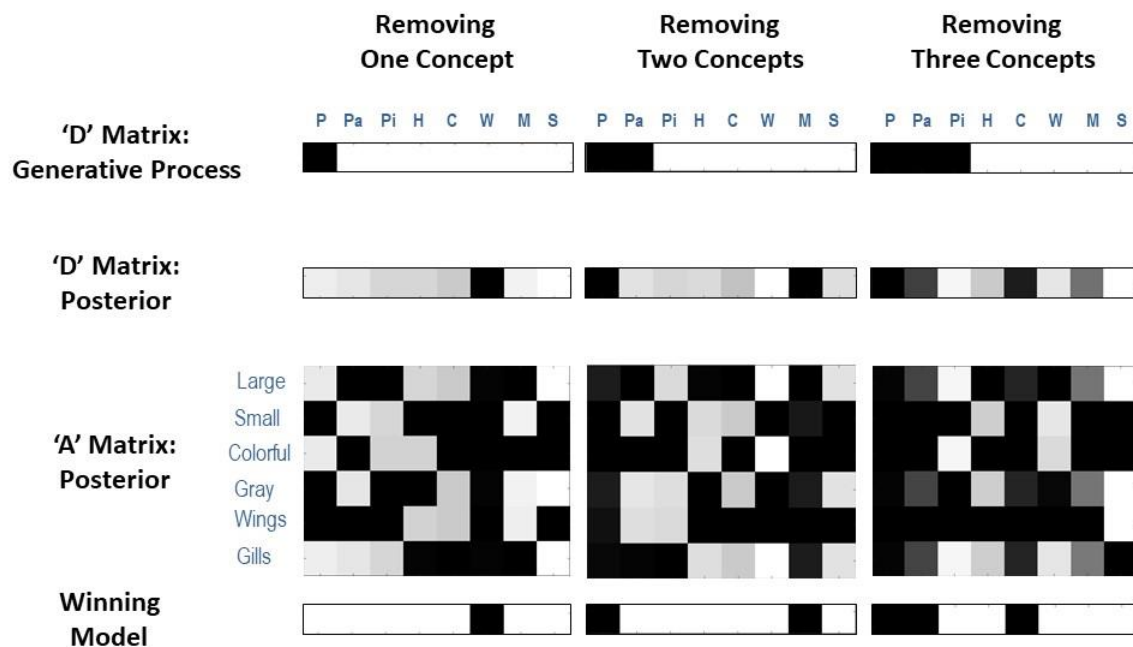
541 unused synaptic connections. In this (biological) setting, the concentration
542 parameters of the implicit Dirichlet distributions can be thought of as synaptic tags.
543 For a technical description of Bayesian model reduction techniques and their
544 proposed neural implementation, see (KJ Friston, Lin, et al., 2017; Hobson & Friston,
545 2012; Hobson, Hong, & Friston, 2014a); also see the left panel of Figure 8 and the
546 associated legend further below for some additional details).

547 The posterior concentration parameters were compared to the prior
548 distribution for a full model (i.e., a flat distribution over 8 concepts) and prior
549 distributions for possible reduced models (i.e., which retained different possible
550 combinations of some but not all concepts; technically, reduced models were
551 defined such that the to-be-eliminated concepts were less likely than the to-be-
552 retained concepts). If Bayesian model reduction provided more evidence for one or
553 more reduced models, the reduced model with the most evidence was selected.
554 Note: an alternative would be to perform model reduction on the 'A' matrix, but this
555 is more complex due to a larger space of possible reduced models; it also does not
556 address the question of the number of hidden state levels to retain in a
557 straightforward manner.

558 In our first simulation, we presented our agent with all animals except for
559 parakeets with equal probability over 2000 trials. When compared to the full model,
560 the winning model corresponded to the correct 7-animal model matching the
561 generative process in 6/8 cases (log evidence differences ranged from -3.12 to -8.3),
562 and in 2/8 cases it instead selected a 6-animal model due to a failure to distinguish
563 between 2 specific concepts during learning (log evidence differences = -5.30, -

564 7.70). Figure 6 illustrates the results of a representative successful case). In the
 565 successful cases, this would correctly license the removal of changes in the model's
 566 'A' and 'D' matrix parameters for the 8th animal concept during learning in the
 567 previous trials. Similar results were obtained whenever any single animal type was
 568 absent from the generative process.

Bayesian Model Reduction After Learning



569

570 Figure 6. Representative illustrations of simulations in which the agent performed Bayesian
 571 model reduction after learning. In these simulations, the agent was first exposed to 2000
 572 trials in which either 7, 6, or 5 animals were actually presented (i.e., illustrated in the top
 573 row, where only the white columns had nonzero probabilities in the generative process). In
 574 each case, model reduction was often successful at identifying the reduced model with the
 575 correct number of animal types presented (bottom row, where black columns should be
 576 removed) based on how much evidence it provided for the posterior distribution over
 577 hidden states learned by the agent (2nd row). This would license the agent to reset
 578 the unneeded columns in its likelihood mapping (3rd row) to their initial state (i.e., a
 579 flat distribution over features) such that it could be engaged if/when a new type of
 580 animal began to be observed (i.e., as in the simulations illustrated in the previous
 581 sections).

582

583 In a second simulation, the generative process contained 2 birds and all 4
584 fish. Here, the correct reduced model was correctly selected in 6/8 simulations (log
585 evidence differences range from -.96 to -8.24, with magnitudes greater than -3 in
586 5/6 cases), whereas it incorrectly selected the 5-animal model in 2 cases (log
587 evidence differences = -3.54, -4.50). In a third simulation, the generative process
588 contained 1 bird and all 4 fish. Here, the correct reduced model had the most
589 evidence in only 3/8 simulations (log evidence differences = -4.10, -4.11, -5.48),
590 whereas a 6-animal model was selected in 3/8 cases and a 3-animal and 7-animal
591 model were each selected once (log evidence differences > -3.0). Figure 6 also
592 illustrates representative examples of correct model recovery in these 2nd and 3rd
593 simulations.

594 While we have used the terms ‘correct’ and ‘incorrect’ above to describe the
595 model used to generate the data, we acknowledge that ‘all models are wrong’ (Box,
596 Hunter, & Hunter, 2005), and that the important question is not whether we can
597 recover the ‘true’ process used to generate the data, but whether we can arrive at
598 the simplest but accurate explanation for these data. The failures to recover the
599 ‘true’ model highlighted above may reflect that a process other than that used to
600 generate the data could have been used to do so in a simpler way. Simpler here
601 means we would have to diverge to a lesser degree, from our prior beliefs, in order
602 to explain the data under a given model, relative to a more complex model. It is
603 worth highlighting the importance of the word *prior* in the previous sentence. This
604 means that the simplicity of the model is sensitive to our prior beliefs about it. To

605 illustrate this, we repeated the same model comparisons as above, but with precise
606 beliefs in an 'A' matrix that complies with that used to generate the data. Specifically,
607 we repeated the three simulations above but only enabled 'D' matrix learning (i.e.,
608 the model was already equipped with the 'A' matrix of the generative process). In
609 each case, Bayesian model reduction now uniquely identified the correct reduced
610 model in 100% of cases.

611 These results demonstrate that – after a naïve model has expanded its hidden
612 state space to include likelihood mappings and initial state priors for a number of
613 concept categories – Bayesian model reduction can subsequently be used to
614 eliminate any parameter updates accrued for *one or two* redundant concept
615 categories. In practice, the 'A' and 'D' concentration parameters for these redundant
616 categories could be reset to their default pre-learning values – and could then be re-
617 engaged if new patterns of observations were repeatedly observed in the future.
618 However, when three concepts should be removed, Bayesian model reduction was
619 much less reliable. This appeared to be due to imperfect 'A' matrix learning, when
620 occurring simultaneously with the (resultingly noisy) accumulation of prior
621 expectations over hidden states – as a fully precise 'A' matrix led to correct model
622 reduction in every case tested (i.e., suggesting that this type of model reduction
623 procedure could be improved by first allowing state-observation learning to
624 proceed alone, then subsequently allowing the model to learn prior expectations
625 over hidden states, which could then be used in model reduction).

626

627 **Can concept acquisition allow for generalization?**

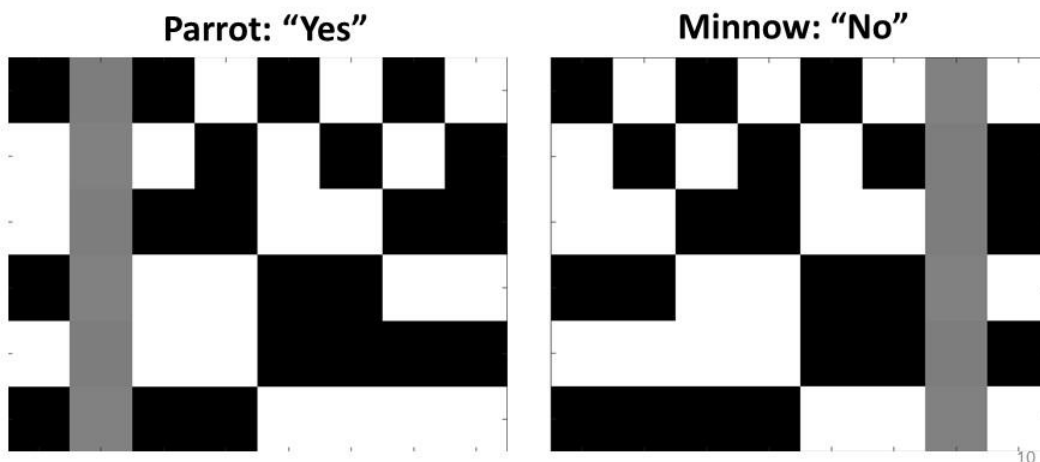
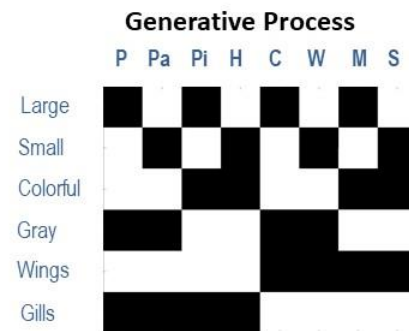
628 One important ability afforded by concept learning is generalization. In a
629 final set of simulations, we asked if our model of concept knowledge could account
630 for generalization. To do so, we altered the model such that it no longer reported
631 what it saw, but instead had to answer a question that depended on generalization
632 from particular cross-category feature combinations. Specifically, the model was
633 shown particular animals and asked: “could this be seen from a distance?” The
634 answer to this question depended on both size and color, such that the answer was
635 yes only for colorful, large animals (i.e., assuming small or gray animals would blend
636 in with the sky or water and be missed).

637 Crucially, this question was asked of animals that the model had not been
638 exposed to, such that it had to generalize from knowledge it already possessed (see
639 Figure 7). To simulate and test for this ability, we equipped the model’s ‘A’ matrix
640 with expert knowledge of 7 out of the 8 animals (i.e., as if these concepts had been
641 learned previously, as in our simulations above). The 8th animal was unknown to the
642 agent, in that it’s likelihood mapping was set such that the 8th animal state “slot”
643 predicted all observations roughly equally. In one variant, the model possessed all
644 concepts except for “parrot,” and it knew that the answer to the question was yes for
645 “whale shark” but not for any other concept it knew. To simulate one-shot
646 generalization, learning was disabled and a parrot (which it had never seen before)
647 was presented 20 times to see if it would correctly generalize and answer “yes” in a
648 reliable manner. In another variant, the model had learned all concepts except
649 “minnow” and was tested the same way to see if it would reliably provide the
650 correct “no” response.

651 Here, we observed that in both of these cases (as well as all others we tested)
652 the model generalized remarkably well. It answered “yes” and “no” correctly in
653 100% of trials. Thus, the agent did not simply possess concepts to explain things it
654 saw. It instead demonstrated generalizable knowledge and could correctly answer
655 questions when seeing a novel stimulus.
656

Generalization:

Could this novel animal
be seen from a distance?



657

658 Figure 7. Depiction of simulations in which we tested the agents ability to generalize from
659 prior knowledge and correctly answered questions about new animals to which it had not
660 previously been exposed. In the simulations, the generative model was modified so that the
661 agent instead chose to report either “yes” or “no” to the question: “could this animal be seen
662 from a distance?” Here, the answer was only yes if the animal was both large and colorful.
663 We observed that when the agent started out with no knowledge of parrots it still correctly
664 answered this question 100% of the time, based only on its knowledge of other animals.
665 Similarly, when it started with no knowledge of minnows, it also correctly reported “no”
666 100% of the time. Thus, the agent was able to generalize from prior knowledge with no
667 additional learning.

668

669

670 **Open questions and relation to other theoretical accounts of concept learning**

671

672 As our simulations show, this model allows for learning novel concepts (i.e.,
673 novel hidden states) based on assigning one or more ‘open slots’ that can be utilised
674 to learn novel feature combinations. In a simple example, we have shown that this
675 setup offers a potential computational mechanism for ‘model expansion’; i.e., the
676 process of expanding a state space to account for novel instances in perceptual
677 categorisation. We also illustrated how this framework can be combined with model
678 reduction, which may be a mechanism for ‘re-setting’ these open slots based on
679 recent experience.

680 This provides a first step towards understanding how agents flexibly expand
681 or reduce their model to adapt to ongoing experience. Yet, several open questions
682 remain, which have partly been addressed in previous work. For example, the
683 proposed framework resonates with previous similarity-based accounts of concept
684 learning. Previous work has proposed a computational framework for arbitrating
685 between assigning an observation to a previously formed memory or forming a
686 novel (hidden) state representation (S. J. Gershman, Monfils, Norman, & Niv, 2017),
687 based on evidence that this observation was sampled from an existing or novel
688 latent state. This process is conceptually similar to our application of Bayesian
689 model reduction over states. In the present framework, concept learning relies on a
690 process based on inference and learning. First, agents have to *infer* whether ongoing

691 observations can be sufficiently explained by existing conceptual knowledge – or
692 speak to the presence of a novel concept that motivates the use of an ‘open slot’.
693 This process is cast as inference on (hidden) states. Second, if the agent infers that
694 there is a novel concept that explains current observations, it has to *learn* about the
695 specific feature configuration of that concept (i.e., novel state). This process
696 highlights the interplay between inference, which allows for the acquisition of
697 knowledge on a relatively short timescale, and learning, which allows for knowledge
698 acquisition on a longer and more stable timescale.

699 Similar considerations apply to the degree of ‘similarity’ of observations. In
700 the framework proposed here, we have assumed that the feature space of
701 observations is already learned and fixed. However, these feature spaces have to be
702 learned in the first place, which implies learning the underlying components or
703 feature dimensions that define observations. This relates closely to notion of
704 structure learning as dimensionality reduction based on covariance between
705 observations, as prominently discussed in the context of spatial navigation (Behrens
706 et al., 2018; Dordek, Soudry, Meir, & Derdikman, 2016; Stachenfeld et al., 2016;
707 Whittington, Muller, Mark, Barry, & Behrens, 2018).

708 Another important issue is how such abstract conceptual knowledge is
709 formed across different contexts or tasks. For example, the abstract concept of a
710 ‘bird’ will be useful for learning about the fauna in a novel environment, but specific
711 types of birds – tied to a previous context – might be less useful in this regard. This
712 speaks to the formation of abstract, task-general knowledge that results from
713 training across different tasks, as recently discussed in the context of meta-

714 reinforcement learning (Ritter, Wang, Kurth-Nelson, & Botvinick, 2018; J X Wang et
715 al., 2016) with a putative link to the prefrontal cortex (Jane X. Wang et al., 2018). In
716 the present framework, such task-general knowledge would speak to the formation
717 of a hierarchical organisation that allows for the formation of conceptual knowledge
718 both within and across contexts. Also note that our proposed framework depends
719 on a pre-defined state space, including a pre-defined set of ‘open slots’ that allow for
720 novel context learning. The contribution of the present framework is to show how
721 these ‘open slots’ can be used for novel concept learning and be re-set based on
722 model reduction. It will be important to extend this approach towards learning the
723 structure of these models in the first place, including the appropriate number of
724 ‘open slots’ (i.e., columns of the A-matrix) for learning in a particular content
725 domain and the relevant feature dimensions of observations (i.e., rows of A-matrix).

726 This corresponds to a potentially powerful and simple application of
727 Bayesian model reduction, in which candidate models (i.e., reduced forms of a full
728 model) are readily identifiable based upon the similarity between the likelihoods
729 conditioned upon different hidden states. If two or more likelihoods are sufficiently
730 similar, the hidden states can be merged (by assigning the concentration
731 parameters accumulated during experience-dependent learning to one or other of
732 the hidden states). The ensuing change in model evidence scores the reduction in
733 complexity. If this reduction is greater than the loss of accuracy – in relation to
734 observations previously encountered – Bayesian model reduction will, effectively,
735 merge one state into another; thereby freeing up a state for the learning of new

736 concepts. We will demonstrate this form of structure learning via Bayesian model
737 reduction in future work.

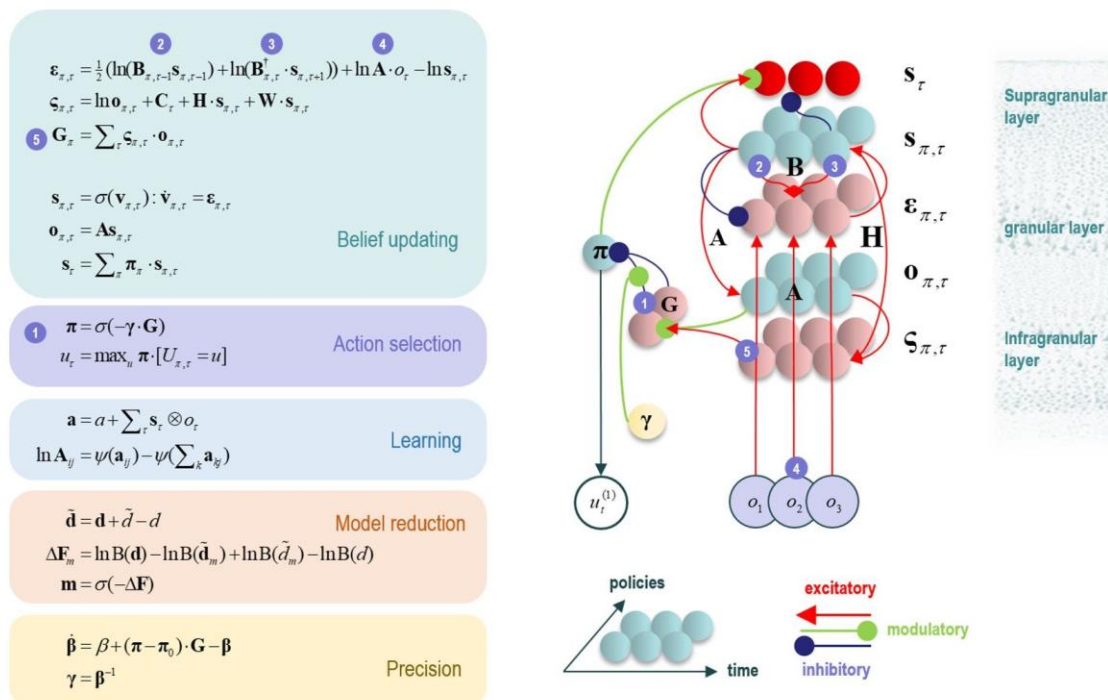
738

739 **Potential advantages of the approach**

740 The present approach may offer some potential theoretical and empirical
741 advantages in comparison to previous work. One theoretical advantage corresponds
742 to the parsimony of casting this type of structure learning as an instance of Bayesian
743 model selection. When integrated with other aspects of the active inference
744 framework, this entails that perceptual inference, active learning, and structure
745 learning are all expressions of the same principle; namely, the minimization of
746 variational free energy, over three distinct timescales. A second, related theoretical
747 advantage is that, when this type of structure learning is cast as Bayesian model
748 selection/reduction, there is no need to invoke additional procedures or schemes
749 (e.g., nonparametric Bayes or ‘stick breaking’ processes; (S. Gershman & Blei,
750 2012)). Instead, a generative model with the capacity to represent a sufficiently
751 complex world will automatically learn causal structure in a way that contextualizes
752 active inference within active learning, and active learning within structure
753 learning.

754 One potential empirical advantage of the present approach stems from the
755 fact that active inference models have a plausible biological basis that affords
756 testable neurobiological predictions. Specifically, these models have well-articulated
757 companion micro-anatomical neural process theories, based on commonly used
758 message-passing algorithms (KJ Friston, FitzGerald, et al., 2017; Parr & Friston,

759 2018; Parr, Markovic, Kiebel, & Friston, 2019). In these process theories, for
 760 example, the activation level of different neural populations (typically portrayed as
 761 consisting of different cortical columns) can encode posterior probability estimates
 762 over different hidden states. These activation levels can then be updated by synaptic
 763 inputs with particular weights that convey the conditional probabilities encoded in
 764 the 'A' and 'B' (among other) matrices described above, where active learning then
 765 corresponds to associative synaptic plasticity. Phasic dopamine responses also play
 766 a particular role in these models, by reporting changes in policy precision (i.e., the
 767 degree of confidence in one policy over others) upon new observations (see Figure 8
 768 and the associated legend for more details).
 769



770

771 Figure 8. This figure illustrates the mathematical framework of active inference and
 772 associated neural process theory used in the simulations described in this paper. The
 773 differential equations in the left panel approximate Bayesian belief updating within the
 774 graphical model depicted in the right panel of Figure 1 via a gradient descent on free energy

775 (F). The right panel also illustrates the proposed neural basis by which neurons making up
776 cortical columns could implement these equations. The equations have been expressed in
777 terms of two types of prediction errors. State prediction errors (ϵ) signal the difference
778 between the (logarithms of) expected states (\mathbf{s}) under each policy and time point—and the
779 corresponding predictions based upon outcomes/observations (\mathbf{A} matrix) and the
780 (preceding and subsequent) hidden states (\mathbf{B} matrix, and, although not written, the \mathbf{D}
781 matrix for the initial hidden states at the first time point). These represent prior and
782 likelihood terms respectively – also marked as messages 2, 3, and 4, which are depicted as
783 being passed between neural populations (colored balls) via particular synaptic
784 connections in the right panel. These (prediction error) signals drive depolarization (\mathbf{v}) in
785 those neurons encoding hidden states (\mathbf{s}), where the probability distribution over hidden
786 states is then obtained via a softmax (normalized exponential) function (σ). Outcome
787 prediction errors (\mathbf{c}) instead signal the difference between the (logarithms of) expected
788 observations (\mathbf{o}) and those predicted under prior preferences (\mathbf{C}). This term additionally
789 considers the expected ambiguity or conditional entropy (\mathbf{H}) between states and outcomes
790 as well as a novelty term (\mathbf{W}) reflecting the degree to which beliefs about how states
791 generate outcomes would change upon observing different possible state-outcome
792 mappings (computed from the \mathbf{A} matrix). This prediction error is weighted by the expected
793 observations to evaluate the expected free energy (\mathbf{G}) for each policy (π), conveyed via
794 message 5. These policy-specific free energies are then integrated to give the policy
795 expectations via a softmax function, conveyed through message 1. Actions at each time
796 point (\mathbf{u}) are then chosen out of the possible actions under each policy (\mathbf{U}) weighted by the
797 value (negative expected free energy) of each policy. In our simulations, the model learned
798 associations between hidden states and observations (\mathbf{A}) via a process in which counts
799 were accumulated (\mathbf{a}) reflecting the number of times the agent observed a particular
800 outcome when she believed that she occupied each possible hidden state. Although not
801 displayed explicitly, learning prior expectations over initial hidden states (\mathbf{D}) is similarly
802 accomplished via accumulation of concentration parameters (\mathbf{d}). These prior expectations
803 reflect counts of how many times the agent believes it previously occupied each possible
804 initial state. Concentration parameters are converted into expected log probabilities using
805 digamma functions (ψ). The way in which Bayesian model reduction was performed in this
806 paper is also written in the lower left (where B indicates a beta function, and \mathbf{m} is the
807 posterior probability of each model). Here, the posterior distribution over initial states (\mathbf{d})
808 is used to assess the difference in the evidence (ΔF) it provides for the number of hidden
809 states in the current model and other possible models characterized by fewer hidden states.
810 Prior concentration parameters are shown in italics, posterior in bold, and those priors and
811 posteriors associated with the reduced model are equipped with a tilde (\sim). As already
812 stated, the right panel illustrates a possible neural implementation of the update equations
813 in the middle panel. In this implementation, probability estimates have been associated
814 with neuronal populations that are arranged to reproduce known intrinsic (within cortical
815 area) connections. Red connections are excitatory, blue connections are inhibitory, and
816 green connections are modulatory (i.e., involve a multiplication or weighting). These
817 connections mediate the message passing associated with the equations in the left panel.
818 Cyan units correspond to expectations about hidden states and (future) outcomes under
819 each policy, while red states indicate their Bayesian model averages (i.e., a “best guess”
820 based on the average of the probability estimates for the states and outcomes across
821 policies, weighted by the probability estimates for their associated policies. Pink units
822 correspond to (state and outcome) prediction errors that are averaged to evaluate expected
823 free energy and subsequent policy expectations (in the lower part of the network). This
824 (neural) network formulation of belief updating means that connection strengths

825 correspond to the parameters of the generative model described in the text. Learning then
826 corresponds to changes in the synaptic connection strengths. Only exemplar connections
827 are shown to avoid visual clutter. Furthermore, we have just shown neuronal populations
828 encoding hidden states under two policies over three time points (i.e., two transitions),
829 whereas in the task described in this paper there are greater number of allowable policies.
830 For more information regarding the mathematics and processes illustrated in this figure,
831 see (KJ Friston, Lin, et al., 2017; KJ Friston, Parr, et al., 2017).

832
833

834

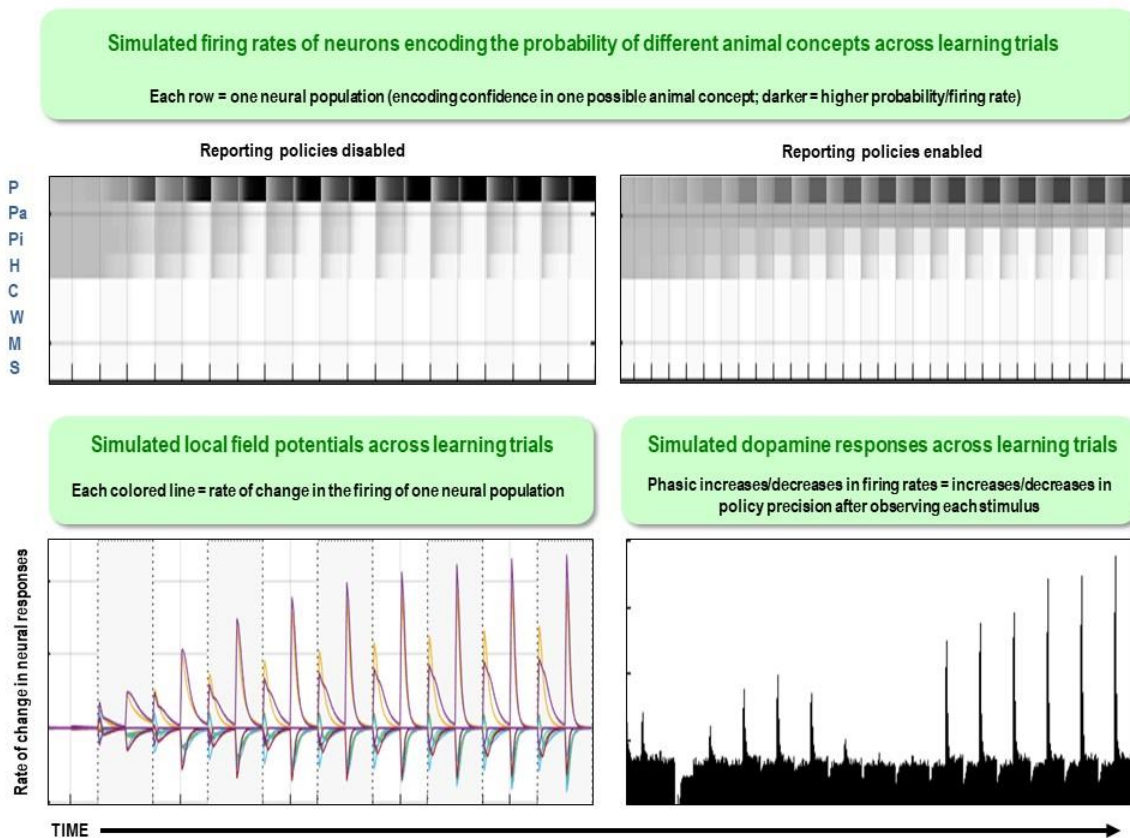
835 Based on these theories, the present model would predict that the brain
836 contains “reserve” cortical columns and synapses (most likely within secondary
837 sensory and association cortices) available to capture new patterns in observed
838 features. To our knowledge, no direct evidence supporting the presence of unused
839 cortical columns in the brain has been observed, although the generation of new
840 neurons (with new synaptic connections) is known to occur in the hippocampus
841 (Chancey et al., 2013). “Silent synapses” have also been observed in the brain, which
842 does appear consistent with this prediction; such synapses can persist into
843 adulthood and only become activated when new learning becomes necessary (e.g.,
844 see (Chancey et al., 2013; Funahashi, Maruyama, Yoshimura, & Komatsu, 2013;
845 Kerchner & Nicoll, 2008)). One way in which this idea of “spare capacity” or
846 “reserve” cortical columns might be tested in the context of neuroimaging would be
847 to examine whether greater levels of neural activation – within conceptual
848 processing regions – are observed after learning additional concepts, which would
849 imply that additional populations of neurons become capable of being activated. In
850 principle, single-cell recording methods might also test for the presence of neurons

851 that remain at baseline firing rates during task conditions, but then become
852 sensitive to new stimuli within the relevant conceptual domain after learning.

853 Figure 9 provides a concrete example of two specific empirical predictions
854 that follow from simulating the neural responses that should be observed within our
855 concept learning task under these process theories. In the left panel, we plot the
856 firing rates (darker = higher firing rate) and local field potentials (rate of change in
857 firing rates) associated with neural populations encoding the probability of the
858 presence of different animals that would be expected across a number of learning
859 trials. In this particular example, the agent began with knowledge of the basic
860 categories of 'bird' and 'fish,' but needed to learn the eight more specific animal
861 categories over 50 interleaved exposures to each animal (only 10 equally spaced
862 learning trials involving the presentation of a parakeet are shown for simplicity). As
863 can be seen, early in learning the firing rates and local field potentials remain at
864 baseline levels; in contrast, as learning progresses, these neural responses take a
865 characteristic shape with more and more positive changes in firing rate in the
866 populations representing the most probable animal, while other populations drop
867 further and further below baseline firing rates.

868 The right panel depicts a similar simulation, but where the agent was
869 allowed to self-report what it saw on each trial (for clarity of illustration, we here
870 show 12 equally spaced learning trials for parakeet over 120 total trials). Enabling
871 policy selection allowed us to simulate expected phasic dopamine responses during
872 the task, corresponding to changes in the precision of the probability distribution
873 over policies after observing a stimulus on each trial. As can be seen, during early

874 trials the model predicts small firing rate increases when the agent is confident in its
875 ability to correctly report the more general animal category after observing a new
876 stimulus, and firing rate decreases when the agent becomes less confident in one
877 policy over others (i.e., as confidence in reporting the specific versus general
878 categories becomes more similar). Larger and larger phasic dopaminergic responses
879 are then expected as the agent becomes more and more confident in its ability to
880 correctly report the specific animal category upon observing a new stimulus. It will
881 be important for future neuroimaging studies to test these predictions in this type of
882 concept learning/stimulus categorization task.
883
884



885

886 Figure 9. Simulated neuronal firing rates, local field potentials, and dopaminergic responses
887 across learning trials based on the neural process theory associated with active inference
888 summarized in Figure 8. The top left panel displays the predicted firing rates (darker =
889 higher firing rate) of neural populations encoding the probability of each hidden state
890 over 50 interleaved exposures to each animal (only 10 equally spaced learning trials
891 involving the presentation of a parakeet are shown for simplicity) in the case where
892 the agent starts out with knowledge of the basic animal categories but must learn
893 the more specific categories. As can be seen, initially each of the four neural
894 populations encoding possible bird categories (i.e., one row per possible category)
895 have equally low firing rates (gray); as learning continues, firing rates increase for
896 the ‘parakeet’ population and decrease for the others. The bottom left panel
897 illustrates the predicted local field potentials (based on the rate of change in firing
898 rates) that would be measured across the task. The top right panel displays the
899 predicted firing rates of neural populations in an analogous simulation in which
900 reporting policies were enabled (for clarity of illustration, we here show 12 equally
901 spaced learning trials for parakeet over 120 total trials). Enabling policy selection
902 allowed us to simulate the phasic dopaminergic responses (reporting changes in the
903 precision of the probability distribution over policies) predicted to occur across
904 learning trials; here the agent first becomes confident in its ability to correctly
905 report the general animal category upon observing a stimulus, then becomes unsure
906 about reporting specific versus general categories, and then becomes confident in
907 its ability to report the specific categories.
908

909

910

Discussion

911

912 The Active Inference formulation of concept learning presented here
913 demonstrates a simple way in which a generative model can acquire both basic and
914 highly granular knowledge of the hidden states/causes in its environment. In
915 comparison to previous theoretical work using active inference (e.g., (M. Mirza,
916 Adams, Mathys, & Friston, 2016; Parr & Friston, 2017; Schwartenbeck, FitzGerald,
917 Mathys, Dolan, & Friston, 2015)), the novel aspect of our model was that it was
918 further equipped with “reserve” hidden states initially devoid of content (i.e., these
919 states started out with uninformative likelihood mappings that predicted all

920 outcomes with roughly equal probability). Over multiple exposures to different
921 stimuli, these hidden states came to acquire conceptual content that captured
922 distinct statistical patterns in the features of those stimuli. This was accomplished
923 via the model's ability to infer when its currently learned hidden states were unable
924 to account for a new observation, leading an unused hidden state column to be
925 engaged that could acquire a new state-observation mapping.

926 Crucially, the model was able to start with some concepts and then expand its
927 representational repertoire to learn others – but would only do so when a new
928 stimulus was observed. This is conceptually similar to nonparametric Bayesian
929 learning models, such as the “Chinese Room” process and the “Indian Buffet”
930 process, that can also infer the need to invoke additional hidden causes with
931 additional data (S. Gershman & Blei, 2012). These statistical learning models do not
932 need to build in additional “category slots” for learning as in our model and can, in
933 principle, entertain infinite state spaces. On the other hand, it is less clear at present
934 how the brain could implement this type of learning. An advantage of our model is
935 that learning depends solely on biologically plausible Hebbian mechanisms (for a
936 possible neural implementation of model reduction, see (KJ Friston, Lin, et al., 2017;
937 Hobson & Friston, 2012; Hobson, Hong, & Friston, 2014b)).

938 The distinction between nonparametric Bayesian learning and the current
939 active learning scheme may be important from a neurodevelopmental perspective
940 as well. In brief, structure learning in this paper starts with a generative model with
941 ‘spare capacity’, where uncommitted or naive conceptual ‘slots’ are used to explain
942 the sensorium, during optimization of free energy or model evidence. In contrast,

943 nonparametric Bayesian approaches add new slots when appropriate. One might
944 imagine that neonates are equipped with brains with ‘spare capacity’ (Baker &
945 Tenenbaum, 2014) that is progressively leveraged during neurodevelopment, much
946 in the spirit of curriculum learning (Al-Muhaideb & Menai, 2011). In this sense, the
947 current approach to structure learning may be better considered as active learning
948 with generative models that are equipped with a large number of hidden states,
949 which are judiciously reduced – via a process of Bayesian model reduction.
950 Furthermore, as in the acquisition of expertise, our model can also begin with broad
951 category knowledge and then subsequently learn finer-grained within-category
952 distinctions, which has received less attention from the perspective of the
953 aforementioned models. Reporting broad versus specific category recognition is
954 also a distinct aspect of our model – driven by differing levels of uncertainty and an
955 expectation (preference) not to incorrectly report a more specific category.

956 Our simulation results also demonstrated that, when combined with
957 Bayesian model reduction, the model can guard against learning too many
958 categories during model expansion – often retaining only the number of hidden
959 causes actually present in its environment – and to keep “reserve” hidden states for
960 learning about new causes if or when they appear. With perfect “expert” knowledge
961 of the possible animal types it could observe (i.e., fully precise likelihood mappings
962 matching the generative process) this was true in general. Interestingly, however,
963 with an imperfectly learned likelihood mapping, model reduction only succeeded
964 when the agent had to remove either 1 or 2 concepts from her model; when 3
965 potential categories needed to be removed, the correct reduced model was

966 identified less than half the time. It would be interesting to empirically test whether
967 similar learning difficulties are present in humans.

968 Neurobiological theories associated with Active Inference also make
969 predictions about the neural basis of this process (Hobson & Friston, 2012; Hobson
970 et al., 2014b). Specifically, during periods of rest (e.g., daydreaming) or sleep, it is
971 suggested that, because sensory information is down-weighted, learning is driven
972 mainly by internal model simulations (e.g., as appears to happen in the phenomenon
973 of hippocampal replay; (Feld & Born, 2017; Lewis, Knoblich, & Poe, 2018; Pfeiffer &
974 Foster, 2013)); this type of learning can accomplish a model reduction process in
975 which redundant model parameters are identified and removed to prevent model
976 over-fitting and promote selection of the most parsimonious model that can
977 successfully account for previous observations. This is consistent with work
978 suggesting that, during sleep, many (but not all) synaptic strength increases
979 acquired in the previous day are attenuated (Tononi & Cirelli, 2014). The role of
980 sleep and daydreaming in keeping “reserve” representational resources available
981 for model expansion could therefore be especially important to concept learning –
982 consistent with the known role of sleep in learning and memory (Ackermann &
983 Rasch, 2014; Feld & Born, 2017; Perogamvros & Schwartz, 2012; Stickgold, Hobson,
984 Fosse, & Fosse, 2001; Walker & Stickgold, 2010).

985 In addition, an emergent feature of our model was its ability to generalize
986 prior knowledge to new stimuli to which it had not previously been exposed. In fact,
987 the model could correctly generalize upon a single exposure to a new stimulus – a
988 type of “one-shot learning” capacity qualitatively similar to that observed in humans

989 (Landau, Smith, & Jones, 1988; E. Markman, 1989; Xu & Tenenbaum, 2007b). While
990 it should be kept in mind that the example we have provided is very simple, it
991 demonstrates the potential usefulness of this novel approach. Some other
992 prominent approaches in machine-learning (e.g., deep learning) tend to require
993 larger amounts of data (Geman et al., 1992; Hinton et al., 2012; LeCun et al., 2015;
994 Lecun et al., 1998; Mnih et al., 2015), and do not learn the rich structure that allows
995 humans to use concept knowledge in a wide variety of generalizable functions
996 (Barsalou, 1983; Biederman, 1987; Feldman, 1997; Jern & Kemp, 2013; A. B.
997 Markman & Makin, 1998; Osherson & Smith, 1981; Ward, 1994; Williams &
998 Lombrozo, 2010). Other recent hierarchical Bayesian approaches in cognitive
999 science have made progress in this domain, however, by modeling concepts as types
1000 of probabilistic programs (Ghahramani, 2015; Goodman, Tenenbaum, &
1001 Gerstenberg, 2015; Lake et al., 2015).

1002 It is important to note that this model is deliberately simple and is meant
1003 only to represent a proof of principle that categorical inference and conceptual
1004 knowledge acquisition can be modeled within this particular neurocomputational
1005 framework. We chose a particular set of feature combinations to illustrate this, but it
1006 remains to be demonstrated that learning in this model would be equally successful
1007 with a larger feature space and set of learnable hidden causes.

1008 Finally, another topic for future work would be the expansion of this type of
1009 model to context-specific learning (e.g., with an additional hidden state factor for
1010 encoding distinct contexts). In such cases, regularities in co-occurring features differ
1011 in different contexts and other cues to context may not be directly observable (e.g.,

1012 the same species of bird could be a slightly different color or size in different parts
1013 of the world that otherwise appear similar) – creating difficulties in inferring when
1014 to update previously learned associations and when to instead acquire competing
1015 associations assigned to new contexts. At present, it is not clear whether the
1016 approach we have illustrated would be successful at performing this additional
1017 function, although the process of inferring the presence of a new hidden state in a
1018 second hidden state factor encoding context would be similar (for related work on
1019 context-dependent contingency learning, see (S. J. Gershman et al., 2017; S.
1020 Gershman, Jones, Norman, Monfils, & Niv, 2013)). Another point worth highlighting
1021 is that we have made particular choices with regard to various model parameters
1022 and the number of observations provided during learning. Further investigations of
1023 the space of these possible parameter settings will be important. With this in mind,
1024 however, our current modelling results could offer additional benefits. For example,
1025 the model’s simplicity could be amenable to empirical studies of saccadic eye
1026 movements toward specific features during novel category learning (e.g. following
1027 the approach of (M. B. Mirza, Adams, Mathys, & Friston, 2018)). This approach could
1028 also be combined with measures of neural activity in humans or other animals,
1029 allowing more direct tests of the predictions highlighted above. In addition, the
1030 introduction of exploratory, novelty-seeking, actions could be used to reduce the
1031 number of samples required for learning, with agents selecting those data that are
1032 most relevant.

1033 In conclusion, the Active Inference scheme we have described illustrates
1034 feature integration in the service of conceptual inference: it can successfully

1035 simulate simple forms of concept acquisition and concept differentiation (i.e.
1036 increasing granularity), and it spontaneously affords one-shot generalization.
1037 Finally, it speaks to empirical work in which behavioral tasks could be designed to
1038 fit such models, which would allow investigation of individual differences in concept
1039 learning and its neural basis. For example, such a model can simulate (neuronal)
1040 belief updating to predict neuroimaging responses as we illustrated above; i.e., to
1041 identify the neural networks engaged in evidence accumulation and learning
1042 (Schwartenbeck et al., 2015). In principle, the model parameters (e.g., ‘A’ matrix
1043 precision) can also be fit to behavioral choices and reaction times – and thereby
1044 phenotype subjects in terms of the priors under which they infer and learn
1045 (Schwartenbeck & Friston, 2016). This approach could therefore advance
1046 neurocomputational approaches to concept learning in several directions.

1047

1048 **Software note**

1049 Although the generative model – specified by the various matrices described in this
1050 paper – changes from application to application, the belief updates are generic and
1051 can be implemented using standard routines (here **spm_MDP_VB_X.m**). These
1052 routines are available as Matlab code in the SPM academic
1053 software: <http://www.fil.ion.ucl.ac.uk/spm/>. The simulations in this paper can be
1054 reproduced (and customised) via running the Matlab code included here is
1055 supplementary material (**Concepts_model.m**).

1056

1057

References

1058

1059

1060 Ackermann, S., & Rasch, B. (2014). Differential Effects of Non-REM and REM Sleep
1061 on Memory Consolidation? *Current Neurology and Neuroscience Reports*, *14*(2),
1062 430. <https://doi.org/10.1007/s11910-013-0430-8>

1063 Al-Muhaideb, S., & Menai, M. E. B. (2011). Evolutionary computation approaches to
1064 the Curriculum Sequencing problem. *Natural Computing*, *10*(2), 891–920.
1065 <https://doi.org/10.1007/s11047-010-9246-5>

1066 Baker, C., & Tenenbaum, J. (2014). Modeling human plan recognition using Bayesian
1067 theory of mind. In G. Sukthankar, C. Geib, H. Dui, D. Pynadath, & R. Goldman
1068 (Eds.), *Plan, activity, and intent recognition* (pp. 177–204). Boston: Morgan
1069 Kaufmann.

1070 Barsalou, L. W. (1983). Ad hoc categories. *Memory & Cognition*, *11*(3), 211–227.
1071 <https://doi.org/10.3758/bf03196968>

1072 Behrens, T. E. J., Muller, T. H., Whittington, J. C. R., Mark, S., Baram, A. B., Stachenfeld,
1073 K. L., & Kurth-Nelson, Z. (2018). What Is a Cognitive Map? Organizing
1074 Knowledge for Flexible Behavior. *Neuron*, *100*(2), 490–509.
1075 <https://doi.org/10.1016/J.NEURON.2018.10.002>

1076 Biederman, I. (1987). Recognition-by-components: a theory of human image
1077 understanding. *Psychological Review*, *94*(2), 115–147.
1078 <https://doi.org/10.1037/0033-295X.94.2.115>

1079 Botvinick, M. M., Niv, Y., & Barto, A. C. (2009). Hierarchically organized behavior and
1080 its neural foundations: A reinforcement learning perspective. *Cognition*, *113*(3),

- 1081 262–280. <https://doi.org/10.1016/J.COGNITION.2008.08.011>
- 1082 Box, G. E., Hunter, J. S., & Hunter, W. G. (2005). *Statistics for experimenters. Wiley*
- 1083 *Series in Probability and Statistics*. Hoboken, NJ.: Wiley.
- 1084 Brown, T. H., Zhao, Y., & Leung, V. (2010). Hebbian plasticity. In *Encyclopedia of*
- 1085 *Neuroscience* (pp. 1049–1056). [https://doi.org/10.1016/B978-008045046-](https://doi.org/10.1016/B978-008045046-9.00796-8)
- 1086 9.00796-8
- 1087 Chancey, J., Adlaf, E., Sapp, M., Pugh, P., Wadiche, J., & Overstreet-Wadiche, L. (2013).
- 1088 GABA depolarization is required for experience-dependent synapse unsilencing
- 1089 in adult-born neurons. *The Journal of Neuroscience : The Official Journal of the*
- 1090 *Society for Neuroscience*, 33(15), 6614–6622.
- 1091 <https://doi.org/10.1523/JNEUROSCI.0781-13.2013>
- 1092 Conant, C., & Ashbey, W. (1970). Every good regulator of a system must be a model
- 1093 of that system. *International Journal of Systems Science*, 1(2), 89–97.
- 1094 <https://doi.org/10.1080/00207727008920220>
- 1095 Cornish, N. J., & Littenberg, T. B. (2007). Tests of Bayesian model selection
- 1096 techniques for gravitational wave astronomy. *Physical Review D*, 76(8), 083006.
- 1097 <https://doi.org/10.1103/PhysRevD.76.083006>
- 1098 Dordek, Y., Soudry, D., Meir, R., & Derdikman, D. (2016). Extracting grid cell
- 1099 characteristics from place cell inputs using non-negative principal component
- 1100 analysis. *ELife*, 5(MARCH2016), 1–36. <https://doi.org/10.7554/eLife.10094>
- 1101 Feld, G., & Born, J. (2017). Sculpting memory during sleep: concurrent consolidation
- 1102 and forgetting. *Current Opinion in Neurobiology*, 44, 20–27.
- 1103 <https://doi.org/10.1016/J.CONB.2017.02.012>

- 1104 Feldman, J. (1997). The Structure of Perceptual Categories. *Journal of Mathematical*
1105 *Psychology*, 41(2), 145–170. <https://doi.org/10.1006/jmps.1997.1154>
- 1106 Friston, K. J., Litvak, V., Oswal, A., Razi, A., Stephan, K. E., van Wijk, B. C. M., ...
1107 Zeidman, P. (2016). Bayesian model reduction and empirical Bayes for group
1108 (DCM) studies. *NeuroImage*, 128, 413–431.
1109 <https://doi.org/10.1016/J.NEUROIMAGE.2015.11.015>
- 1110 Friston, Karl, Mattout, J., Trujillo-Barreto, N., Ashburner, J., & Penny, W. (2007).
1111 Variational free energy and the Laplace approximation. *NeuroImage*, 34(1),
1112 220–234. <https://doi.org/10.1016/J.NEUROIMAGE.2006.08.035>
- 1113 Friston, Karl, Parr, T., & Zeidman, P. (2018). Bayesian model reduction. Retrieved
1114 from <http://arxiv.org/abs/1805.07092>
- 1115 Friston, Karl, & Penny, W. (2011). Post hoc Bayesian model selection. *NeuroImage*,
1116 56(4), 2089–2099. <https://doi.org/10.1016/J.NEUROIMAGE.2011.03.062>
- 1117 Friston, KJ. (2010). The free-energy principle: a unified brain theory? *Nature*
1118 *Reviews. Neuroscience*, 11(2), 127–138. <https://doi.org/10.1038/nrn2787>
- 1119 Friston, KJ, FitzGerald, T., Rigoli, F., Schwartenbeck, P., O Doherty, J., & Pezzulo, G.
1120 (2016). Active inference and learning. *Neuroscience and Biobehavioral Reviews*,
1121 68, 862–879. <https://doi.org/10.1016/j.neubiorev.2016.06.022>
- 1122 Friston, KJ, FitzGerald, T., Rigoli, F., Schwartenbeck, P., & Pezzulo, G. (2017). Active
1123 Inference: A Process Theory. *Neural Computation*, 29(1), 1–49.
1124 https://doi.org/10.1162/NECO_a_00912
- 1125 Friston, KJ, Lin, M., Frith, C., Pezzulo, G., Hobson, J., & Ondobaka, S. (2017). Active
1126 Inference, Curiosity and Insight. *Neural Computation*, 29(10), 2633–2683.

- 1127 https://doi.org/10.1162/neco_a_00999
- 1128 Friston, KJ, Parr, T., & de Vries, B. (2017). The graphical brain: Belief propagation
1129 and active inference. *Network Neuroscience*, 1(4), 381–414.
- 1130 https://doi.org/10.1162/NETN_a_00018
- 1131 Funahashi, R., Maruyama, T., Yoshimura, Y., & Komatsu, Y. (2013). Silent synapses
1132 persist into adulthood in layer 2/3 pyramidal neurons of visual cortex in dark-
1133 reared mice. *Journal of Neurophysiology*, 109(8), 2064–2076.
- 1134 <https://doi.org/10.1152/jn.00912.2012>
- 1135 Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural Networks and the
1136 Bias/Variance Dilemma. *Neural Computation*, 4(1), 1–58.
- 1137 <https://doi.org/10.1162/neco.1992.4.1.1>
- 1138 Gershman, S., & Blei, D. (2012). A tutorial on Bayesian nonparametric models.
1139 *Journal of Mathematical Psychology*, 56(1), 1–12.
- 1140 <https://doi.org/10.1016/J.JMP.2011.08.004>
- 1141 Gershman, S. J., Monfils, M.-H., Norman, K. A., & Niv, Y. (2017). The computational
1142 nature of memory modification. *ELife*, 6. <https://doi.org/10.7554/eLife.23763>
- 1143 Gershman, S. J., & Niv, Y. (2010). Learning latent structure: carving nature at its
1144 joints. *Current Opinion in Neurobiology*, 20(2), 251–256.
- 1145 <https://doi.org/10.1016/J.CONB.2010.02.008>
- 1146 Gershman, S., Jones, C., Norman, K., Monfils, M., & Niv, Y. (2013). Gradual extinction
1147 prevents the return of fear: implications for the discovery of state. *Frontiers in*
1148 *Behavioral Neuroscience*, 7, 164. <https://doi.org/10.3389/fnbeh.2013.00164>
- 1149 Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence.

- 1150 *Nature*, 521(7553), 452–459. <https://doi.org/10.1038/nature14541>
- 1151 Goodman, N. D., Tenenbaum, J. B., & Gerstenberg, T. (2015). *Concepts: New*
- 1152 *Directions*. (E. Margolis & S. Laurence, Eds.). Cambridge, MA: MIT Press.
- 1153 Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A., Jaitly, N., ... Kingsbury, B. (2012).
- 1154 Deep Neural Networks for Acoustic Modeling in Speech Recognition: The
- 1155 Shared Views of Four Research Groups. *IEEE Signal Processing Magazine*, 29(6),
- 1156 82–97. <https://doi.org/10.1109/MSP.2012.2205597>
- 1157 Hobson, J., & Friston, K. (2012). Waking and dreaming consciousness:
- 1158 neurobiological and functional considerations. *Progress in Neurobiology*, 98(1),
- 1159 82–98. <https://doi.org/10.1016/j.pneurobio.2012.05.003>
- 1160 Hobson, J., Hong, C.-H., & Friston, K. (2014a). Virtual reality and consciousness
- 1161 inference in dreaming. *Frontiers in Psychology*, 5, 1133.
- 1162 <https://doi.org/10.3389/fpsyg.2014.01133>
- 1163 Hobson, J., Hong, C.-H., & Friston, K. (2014b). Virtual reality and consciousness
- 1164 inference in dreaming. *Frontiers in Psychology*, 5, 1133.
- 1165 <https://doi.org/10.3389/fpsyg.2014.01133>
- 1166 Jern, A., & Kemp, C. (2013). A probabilistic account of exemplar and category
- 1167 generation. *Cognitive Psychology*, 66(1), 85–125.
- 1168 <https://doi.org/10.1016/j.cogpsych.2012.09.003>
- 1169 Kemp, C., Perfors, A., & Tenenbaum, J. (2007). Learning overhypotheses with
- 1170 hierarchical Bayesian models. *Developmental Science*, 10(3), 307–321.
- 1171 <https://doi.org/10.1111/j.1467-7687.2007.00585.x>
- 1172 Kerchner, G., & Nicoll, R. (2008). Silent synapses and the emergence of a

- 1173 postsynaptic mechanism for LTP. *Nature Reviews. Neuroscience*, 9(11), 813–
1174 825. <https://doi.org/10.1038/nrn2501>
- 1175 Lake, B., Salakhutdinov, R., & Tenenbaum, J. (2015). Human-level concept learning
1176 through probabilistic program induction. *Science (New York, N.Y.)*, 350(6266),
1177 1332–1338. <https://doi.org/10.1126/science.aab3050>
- 1178 Landau, B., Smith, L. B., & Jones, S. S. (1988). The importance of shape in early lexical
1179 learning. *Cognitive Development*, 3(3), 299–321.
1180 [https://doi.org/10.1016/0885-2014\(88\)90014-7](https://doi.org/10.1016/0885-2014(88)90014-7)
- 1181 LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–
1182 444. <https://doi.org/10.1038/nature14539>
- 1183 Lecun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning
1184 applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278–2324.
1185 <https://doi.org/10.1109/5.726791>
- 1186 Lewis, P., Knoblich, G., & Poe, G. (2018). How Memory Replay in Sleep Boosts
1187 Creative Problem-Solving. *Trends in Cognitive Sciences*, 22(6), 491–503.
1188 <https://doi.org/10.1016/j.tics.2018.03.009>
- 1189 MacKay, D. J. C., & Peto, L. C. B. (1995). A hierarchical Dirichlet language model.
1190 *Natural Language Engineering*, 1(03), 289–308.
1191 <https://doi.org/10.1017/S1351324900000218>
- 1192 Markman, A. B., & Makin, V. S. (1998). Referential communication and category
1193 acquisition. *Journal of Experimental Psychology. General*, 127(4), 331–354.
1194 <https://doi.org/10.1037/0096-3445.127.4.331>
- 1195 Markman, E. (1989). *Categorization and Naming in Children*. Cambridge, MA: MIT

- 1196 Press.
- 1197 McKay, R., & Dennett, D. (2009). The evolution of misbelief. *Behavioral and Brain*
1198 *Sciences*, 32(06), 493. <https://doi.org/10.1017/S0140525X09990975>
- 1199 McNicholas, P. (2016). Model-Based Clustering. *Journal of Classification*, 33(3), 331–
1200 373. <https://doi.org/10.1007/s00357-016-9211-9>
- 1201 Mirza, M., Adams, R., Mathys, C., & Friston, K. (2016). Scene Construction, Visual
1202 Foraging, and Active Inference. *Frontiers in Computational Neuroscience*, 10, 56.
1203 <https://doi.org/10.3389/fncom.2016.00056>
- 1204 Mirza, M. B., Adams, R. A., Mathys, C., & Friston, K. J. (2018). Human visual
1205 exploration reduces uncertainty about the sensed world. *PLOS ONE*, 13(1),
1206 e0190429. <https://doi.org/10.1371/journal.pone.0190429>
- 1207 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ...
1208 Hassabis, D. (2015). Human-level control through deep reinforcement learning.
1209 *Nature*, 518(7540), 529–533. <https://doi.org/10.1038/nature14236>
- 1210 Osherson, D. N., & Smith, E. E. (1981). On the adequacy of prototype theory as a
1211 theory of concepts. *Cognition*, 9(1), 35–58. [https://doi.org/10.1016/0010-](https://doi.org/10.1016/0010-0277(81)90013-5)
1212 [0277\(81\)90013-5](https://doi.org/10.1016/0010-0277(81)90013-5)
- 1213 Parr, T., & Friston, K. (2017). Working memory, attention, and salience in active
1214 inference. *Scientific Reports*, 7(1), 14678. [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-017-15249-0)
1215 [017-15249-0](https://doi.org/10.1038/s41598-017-15249-0)
- 1216 Parr, T., & Friston, K. (2018). The Anatomy of Inference: Generative Models and
1217 Brain Structure. *Frontiers in Computational Neuroscience*, 12, 90.
1218 <https://doi.org/10.3389/fncom.2018.00090>

- 1219 Parr, T., Markovic, D., Kiebel, S., & Friston, K. (2019). Neuronal message passing
1220 using Mean-field, Bethe, and Marginal approximations. *Scientific Reports*, 9(1),
1221 1889. <https://doi.org/10.1038/s41598-018-38246-3>
- 1222 Perfors, A., Tenenbaum, J., Griffiths, T., & Xu, F. (2011). A tutorial introduction to
1223 Bayesian models of cognitive development. *Cognition*, 120(3), 302–321.
1224 <https://doi.org/10.1016/j.cognition.2010.11.015>
- 1225 Perogamvros, L., & Schwartz, S. (2012). The roles of the reward system in sleep and
1226 dreaming. *Neuroscience & Biobehavioral Reviews*, 36(8), 1934–1951.
1227 <https://doi.org/10.1016/J.NEUBIOREV.2012.05.010>
- 1228 Pfeiffer, B., & Foster, D. (2013). Hippocampal place-cell sequences depict future
1229 paths to remembered goals. *Nature*, 497(7447), 74–79.
1230 <https://doi.org/10.1038/nature12112>
- 1231 Ritter, S., Wang, J., Kurth-Nelson, Z., & Botvinick, M. (2018). Episodic Control as
1232 Meta-Reinforcement Learning. *BioRxiv*, 360537.
1233 <https://doi.org/10.1101/360537>
- 1234 Salakhutdinov, R., Tenenbaum, J. B., & Torralba, A. (2013). Learning with
1235 hierarchical-deep models. *IEEE Transactions on Pattern Analysis and Machine*
1236 *Intelligence*, 35(8), 1958–1971. <https://doi.org/10.1109/TPAMI.2012.269>
- 1237 Schmidhuber, J. (2006). Developmental robotics, optimal artificial curiosity,
1238 creativity, music, and the fine arts. *Connection Science*, 18(2), 173–187.
1239 <https://doi.org/10.1080/09540090600768658>
- 1240 Schwartenbeck, P., FitzGerald, T., Mathys, C., Dolan, R., & Friston, K. (2015). The
1241 Dopaminergic Midbrain Encodes the Expected Certainty about Desired

- 1242 Outcomes. *Cerebral Cortex*, 25(10), 3434–3445.
- 1243 <https://doi.org/10.1093/cercor/bhu159>
- 1244 Schwartenbeck, P., & Friston, K. (2016). Computational Phenotyping in Psychiatry: A
1245 Worked Example. *ENeuro*, 3(4), ENEURO.0049-0016.2016.
- 1246 <https://doi.org/10.1523/ENEURO.0049-16.2016>
- 1247 Sharot, T. (2011). The optimism bias. *Current Biology*, 21(23), R941–R945.
- 1248 <https://doi.org/10.1016/J.CUB.2011.10.030>
- 1249 Stachenfeld, K. L., Botvinick, M. M., & Gershman, S. J. (2016). The hippocampus as a
1250 predictive map.
- 1251 Stickgold, R., Hobson, J., Fosse, R., & Fosse, M. (2001). Sleep, learning, and dreams:
1252 off-line memory reprocessing. *Science*, 294(5544), 1052–1057.
- 1253 <https://doi.org/10.1126/science.1063530>
- 1254 Tervo, D. G. R., Tenenbaum, J. B., & Gershman, S. J. (2016). Toward the neural
1255 implementation of structure learning. *Current Opinion in Neurobiology*, 37, 99–
1256 105. <https://doi.org/10.1016/J.CONB.2016.01.014>
- 1257 Tononi, G., & Cirelli, C. (2014). Sleep and the Price of Plasticity: From Synaptic and
1258 Cellular Homeostasis to Memory Consolidation and Integration. *Neuron*, 81(1),
1259 12–34. <https://doi.org/10.1016/J.NEURON.2013.12.025>
- 1260 Walker, M., & Stickgold, R. (2010). Overnight alchemy: sleep-dependent memory
1261 evolution. *Nature Reviews. Neuroscience*, 11(3), 218; author reply 218.
- 1262 <https://doi.org/10.1038/nrn2762-c1>
- 1263 Wang, J X, Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., ...
1264 Botvinick, M. (2016). Learning to reinforcement learn. *ArXiv:1611.05763*.

- 1265 Wang, Jane X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z., ...
1266 Botvinick, M. (2018). Prefrontal cortex as a meta-reinforcement learning
1267 system. *Nature Neuroscience*, *21*(6), 860–868.
1268 <https://doi.org/10.1038/s41593-018-0147-8>
- 1269 Ward, T. B. (1994). Structured Imagination: the Role of Category Structure in
1270 Exemplar Generation. *Cognitive Psychology*, *27*(1), 1–40.
1271 <https://doi.org/10.1006/cogp.1994.1010>
- 1272 Whittington, J. C. R., Muller, T. H., Mark, S., Barry, C., & Behrens, T. E. J. (2018).
1273 Generalisation of structural knowledge in the hippocampal-entorhinal system.
- 1274 Williams, J. J., & Lombrozo, T. (2010). The role of explanation in discovery and
1275 generalization: evidence from category learning. *Cognitive Science*, *34*(5), 776–
1276 806. <https://doi.org/10.1111/j.1551-6709.2010.01113.x>
- 1277 Wipf, D. P., & Rao, B. D. (2007). An Empirical Bayesian Strategy for Solving the
1278 Simultaneous Sparse Approximation Problem. *IEEE Transactions on Signal*
1279 *Processing*, *55*(7), 3704–3716. <https://doi.org/10.1109/TSP.2007.894265>
- 1280 Xu, F., & Tenenbaum, J. (2007a). Sensitivity to sampling in Bayesian word learning.
1281 *Developmental Science*, *10*(3), 288–297. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-7687.2007.00590.x)
1282 [7687.2007.00590.x](https://doi.org/10.1111/j.1467-7687.2007.00590.x)
- 1283 Xu, F., & Tenenbaum, J. (2007b). Word Learning as Bayesian Inference. *Psychological*
1284 *Review*, *114*(2), 245–272.
- 1285
- 1286