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

2    **as an example case**

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

4

5    [1]Laureate Institute for Brain Research, Tulsa, OK, USA

6    [2]Wellcome Centre for Human Neuroimaging, Institute of Neurology, University

7    College London, WC1N 3BG, UK

8

9

10

11

12    <u>Corresponding Author Information:</u>

13

14    Ryan Smith

15    Laureate Institute for Brain Research

16    6655 S Yale Ave, Tulsa, OK 74136, USA

17    Email: rsmith@laureateinstitute.org

18                                     **Abstract**

19    Within computational neuroscience, the algorithmic and neural basis of structure

20    learning remains poorly understood. Concept learning is one primary example,

21    which requires both a type of internal model expansion process (adding novel

22    hidden states that explain new observations), and a model reduction process

23    (merging different states into one underlying cause and thus reducing model

24    complexity via meta-learning). Although various algorithmic models of concept

25    learning have been proposed within machine learning and cognitive science, many

26    are limited to various degrees by an inability to generalize, the need for very large

27    amounts of training data, and/or insufficiently established biological plausibility.

28    Using concept learning as an example case, we introduce a novel approach for

29    modeling structure learning within the active inference framework and its

30    accompanying neural process theory. This approach is based on the idea that a

31    generative model can be equipped with extra (hidden state or cause) 'slots' that can

32    be engaged when an agent learns about novel concepts. This can be combined with a

33    Bayesian model reduction process, in which any concept learning – associated with

34    these slots – can be reset in favor of a simpler model with higher model evidence.

35    We use simulations to illustrate this model's ability to add new concepts to its state

36    space (with relatively few observations) and increase the granularity of the

37    concepts it currently possesses. We also simulate the predicted neural basis of these

38    processes. We further show that it accomplishes a simple form of 'one-shot'

39    generalization to new stimuli. Although deliberately simple, these results suggest

40    that this general approach to modeling concept learning within active inference

41    research may also offer useful resources in developing neurocomputational models

42    of structure learning more generally.

43    *Keywords*: Model Expansion; Structure Learning; Concepts; Computational

44    Neuroscience; Active Inference

45

## Introduction

46

47       The ability to learn the latent structure of one's environment – such as

48 inferring the existence of hidden causes of regularly observed patterns in co-

49 occurring feature observations – is central to human cognition. For example, we do

50 not simply observe particular sets of colors, textures, shapes, and sizes – we also

51 observe *identifiable object*s such as, say, a 'screwdriver'. If we were tool experts, we

52 might also recognize particular types of screwdrivers (e.g., flat vs. Phillip's head),

53 designed for a particular use. This ability to learn latent structure, such as learning

54 to recognize co-occurring features under conceptual categories (as opposed to just

55 perceiving sensory qualities; e.g., red, round, etc.), is also highly adaptive. Only if we

56 knew an object was a screwdriver could we efficiently infer that it affords putting

57 certain structures together and taking them apart; and only if we knew the specific

58 type of screwdriver could we efficiently infer, say, the artefacts to use it on. Many

59 concepts of this sort require experience-dependent acquisition (i.e., learning).

60       From a computational perspective, the ability to acquire a new concept can

61 be seen as a type of structure learning involving Bayesian model comparison

62 (Botvinick, Niv, & Barto, 2009; S. J. Gershman & Niv, 2010; MacKay & Peto, 1995;

63 Salakhutdinov, Tenenbaum, & Torralba, 2013; Tervo, Tenenbaum, & Gershman,

64 2016). Specifically, concept acquisition can be cast as an agent learning (or

65 inferring) that a new hypothesis (referred to here as a hidden cause or state) should

66 be added to the internal or generative model with which she explains her

67 environment, because existing causes cannot account for new observations (e.g., an

68 agent might start out believing that the only tools are hammers and screwdrivers,

69    but later learn that there are also wrenches). In other words, the structure of the

70    space of hidden causes itself needs to expand to accommodate new patterns of

71    observations. This model expansion process is complementary to a process called

72    Bayesian model reduction (Karl Friston & Penny, 2011); in which the agent can infer

73    that there is redundancy in her model, and a model with fewer states or parameters

74    provides a more parsimonious (i.e. simpler) explanation of observations (KJ Friston,

75    Lin, et al., 2017; Schmidhuber, 2006). For example, in some instances it may be

76    more appropriate to differentiate between fish and birds as opposed to salmon,

77    peacocks and pigeons. This reflects a reduction in model complexity based on a

78    particular feature space underlying observations and thus resonates with other

79    accounts of concept learning as dimensionality reduction (Behrens et al., 2018;

80    Stachenfeld, Botvinick, & Gershman, 2016) – a topic we discuss further below.

81          A growing body of work in a number of domains has approached this

82    problem from different angles. In developmental psychology and cognitive science,

83    for example, probability theoretic (Bayesian) models have been proposed to account

84    for word learning in children and the remarkable human ability to generalize from

85    very few (or even one) examples in which both broader and narrower categorical

86    referents could be inferred (Kemp, Perfors, & Tenenbaum, 2007; Lake,

87    Salakhutdinov, & Tenenbaum, 2015; Perfors, Tenenbaum, Griffiths, & Xu, 2011; Xu &

88    Tenenbaum, 2007a, 2007b). In statistics, a number of nonparametric Bayesian

89    models, such as the "Chinese Room" process and the "Indian Buffet" process, have

90    been used to infer the need for model expansion (S. Gershman & Blei, 2012). There

91    are also related approaches in machine learning, as applied to things like Gaussian

92    mixture models (McNicholas, 2016).

93         Such approaches employ various clustering algorithms, which take sets of

94    data points in a multidimensional space and divide them into separable clusters

95    (e.g., see (Anderson, 1991; Love, Medin, & Gureckis, 2004; Sanborn, Griffiths, &

96    Navarro, 2010)). While many of these approaches assume the number of clusters is

97    known in advance, various goodness-of-fit criteria may be used to determine the

98    optimal number. However, a number of approaches require much larger amounts of

99    data than humans do to learn new concepts (Geman, Bienenstock, & Doursat, 1992;

100   Hinton et al., 2012; LeCun, Bengio, & Hinton, 2015; Lecun, Bottou, Bengio, & Haffner,

101   1998; Mnih et al., 2015). Many also require corrective feedback to learn and yet fail

102   to acquire sufficiently rich conceptual structure to allow for generalization

103   (Barsalou, 1983; Biederman, 1987; Feldman, 1997; Jern & Kemp, 2013; A. B.

104   Markman & Makin, 1998; Osherson & Smith, 1981; Ward, 1994; Williams &

105   Lombrozo, 2010).

106        Approaches to formally modeling structure learning, including concept

107   learning, have not yet been examined within the emerging field of research on

108   Active Inference models within computational neuroscience (KJ Friston, 2010; KJ

109   Friston et al., 2016; KJ Friston, Lin, et al., 2017; KJ Friston, Parr, & de Vries, 2017).

110   This represents one potentially fruitful research avenue that has not yet been

111   examined and, as discussed below, may offer unique advantages in research focused

112   on understanding the neural basis of learning latent structure. In this paper, we

113   explore the potential of this approach. In brief, we conclude that structure learning

114    is an emergent property of active inference (and learning) under generative models

115    with 'spare capacity'; where spare or uncommitted capacity is used to expand the

116    repertoire of representations (Baker & Tenenbaum, 2014), while Bayesian model

117    reduction (KJ Friston, Lin, et al., 2017; Hobson & Friston, 2012) promotes

118    generalization by minimizing model complexity – and releasing representations to

119    replenish 'spare capacity'.

120        From a machine learning perspective, Bayesian model reduction affords the

121    opportunity to consider generative models with a large number of hidden states or

122    latent factors and then optimize the number (or indeed partitions) of latent factors

123    with respect to a variational bound on model evidence. This could be regarded as a

124    bounded form of nonparametric Bayes, in which a potentially infinite number of

125    latent factors are considered; with appropriate (e.g., Indian buffet process) priors

126    over the number of hidden states generating data features[1]. The Bayesian model

127    reduction approach here is bounded in the sense that an upper bound on the

128    number of hidden states is specified prior to structure learning. Furthermore, in

129    virtue of the (biologically plausible) variational schemes used for model reduction,

130    there is no need to explicitly compute model evidence; thereby emulating the

131    computational efficiency of nonparametric Bayes (S. Gershman & Blei, 2012), while

132    accommodating any prior over models.

133        In what follows, we first provide a brief overview of active inference. We then

134    introduce a model of concept learning (using basic and subordinate level animal

---

[1] Generally motivated by starting with a finite parametric model and taking the limit as the number of latent states with more parameters tends to infinity.

135    categories), as a representative example of structure learning. We specifically model

136    cognitive (semantic) processes that add new concepts to a state space and that

137    optimize the granularity of an existing state space. We then establish the validity of

138    this model using numerical analyses of concept learning, when repeatedly

139    presenting a synthetic agent with different animals characterized by different

140    combinations of observable features. We demonstrate how particular approaches

141    combining Bayesian model reduction and expansion can reproduce successful

142    concept learning without the need for corrective feedback – and allow for

143    generalization. We further demonstrate the ability of this model to generate

144    predictions about measurable neural responses based on the neural process theory

145    that accompanies active inference. We conclude with a brief discussion of the

146    implications of this work. Our goal is to present an introductory proof of concept –

147    that could be used as the foundation of future research on the neurocomputational

148    basis of structure learning.

149

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

151

152    **A primer on Active Inference**

153

154            Active Inference suggests that the brain is an inference machine that

155    approximates optimal probabilistic (Bayesian) belief updating across perceptual,

156    cognitive, and motor domains. Active Inference more specifically postulates that the

157    brain embodies an internal model of the world that is "generative" in the sense that

158     it can simulate the sensory data that it should receive if its model of the world is

159     correct. These simulated (predicted) sensory data can be compared to actual

160     observations, and differences between predicted and observed sensations can be

161     used to update the model. Over short timescales (e.g., a single observation) this

162     updating corresponds to inference (perception), whereas on longer timescales it

163     corresponds to learning (i.e., updating expectations about what will be observed

164     later). Another way of putting this is that perception optimizes beliefs about the

165     current state of the world, while learning optimizes beliefs about the relationships

166     between the variables that constitute the world. These processes can be seen as

167     ensuring the generative model (entailed by recognition processes in the brain)

168     remains an accurate model of the world that it seeks to regulate (Conant & Ashbey,

169     1970).

170             Active Inference casts decision-making in similar terms. Actions can be

171     chosen to resolve uncertainty about variables within a generative model (i.e.,

172     sampling from domains in which the model does not make precise predictions),

173     which can prevent anticipated deviations from predicted outcomes. In addition,

174     some expectations are treated as a fixed phenotype of an organism. For example, if

175     an organism did not continue to "expect" to observe certain amounts of food, water,

176     and shelter, then it would quickly cease to exist (McKay & Dennett, 2009) – as it

177     would not pursue those behaviors that fulfill these expectations (c.f. the 'optimism

178     bias' (Sharot, 2011)). Thus, a creature should continually seek out observations that

179     support – or are internally consistent with – its own continued existence. Decision-

180     making can therefore be cast as a process in which the brain infers the sets of

181    actions (policies) that would lead to observations consistent with its own survival-

182    related expectations (i.e., its "prior preferences"). Mathematically, this can be

183    described as selecting sequences of actions (policies) that maximize "Bayesian

184    model evidence" expected under a policy, where model evidence is the (marginal)

185    likelihood that particular sensory inputs would be observed under a given model.

186         In real-world settings, directly computing Bayesian model evidence is

187    generally intractable. Thus, some approximation is necessary. Active Inference

188    proposes that the brain computes a quantity called "variational free energy" that

189    provides a bound on model evidence, such that minimization of free energy

190    indirectly maximizes model evidence (this is exactly the same functional used in

191    machine learning where it is known as an evidence lower bound or ELBO). In this

192    case, decision-making will be approximately (Bayes) optimal if it infers (and enacts)

193    the policy that will minimize expected free energy (i.e., free energy with respect to a

194    policy, where one takes expected future observations into account). Technically,

195    expected free energy is the average free energy under the posterior predictive

196    density over policy-specific outcomes.

197         Expected free energy can be decomposed in different ways that reflect

198    uncertainty and prior preferences, respectively (e.g., epistemic and instrumental

199    affordance or ambiguity and risk). This formulation means that any agent that

200    minimizes expected free energy engages initially in exploratory behavior to

201    minimize uncertainty in a new environment. Once uncertainty is resolved, the agent

202    then exploits that environment to fulfil its prior preferences. The formal basis for

203    Active Inference has been thoroughly detailed elsewhere (KJ Friston, FitzGerald,

204    Rigoli, Schwartenbeck, & Pezzulo, 2017), and the reader is referred there for a full

205    mathematical treatment.

206        When the generative model is formulated as a partially observable Markov

207    decision process (a mathematical framework for modeling decision-making in cases

208    where some outcomes are under the control of the agent and others are not, and

209    where states of the world are not directly known but must be inferred from

210    observations), active inference takes a particular form. Here, the generative model is

211    specified by writing down plausible or allowable policies, hidden states of the world

212    (that must be inferred from observations), and observable outcomes, as well as a

213    number of matrices that define the probabilistic relationships between these

214    quantities (see right panel of figure1).
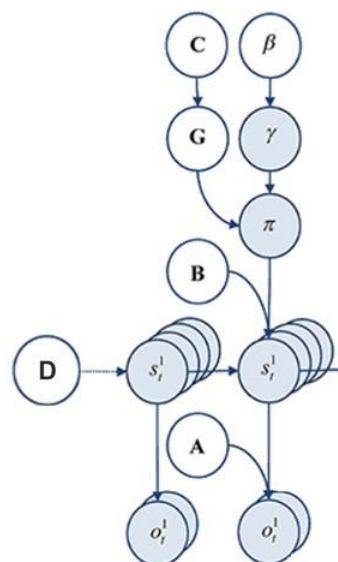


215

216  Figure 1. Left: Illustration of the trial structure performed by the agent. At the first time
217  point, the agent is exposed to one of 8 possible animals that are each characterized by a
218  unique combination of visual features. At the 2nd time point, the agent would then report
219  which animal concept matched that feature combination. The agent could report a specific
220  category (e.g., pigeon, hawk, minnow, etc.) or a general category (i.e., bird or fish) if
221  insufficiently certain about the specific category. See the main text for more details. Right:
222  Illustration of the Markov decision process formulation of active inference used in the
223  simulations described in this paper. The generative model is here depicted graphically, such
224  that arrows indicate dependencies between variables. Here observations ($\mathbf{o}$) depend on
225  hidden states ($\mathbf{s}$), as specified by the $\mathbf{A}$ matrix, and those states depend on both previous
226  states (as specified by the $\mathbf{B}$ matrix,  or the initial states specified by the $\mathbf{D}$ matrix) and the
227  policies ($\boldsymbol{\pi}$) selected by the agent. The probability of selecting a particular policy in turn
228  depends on the expected free energy ($\mathbf{G}$) of each policy with respect to the prior preferences
229  ($\mathbf{C}$) of the agent. The degree to which expected free energy influences policy selection is also
230  modulated by a prior policy precision parameter ($\boldsymbol{\gamma}$), which is in turn dependent on beta ($\boldsymbol{\beta}$)
231  –where higher values of beta promote more randomness in policy selection (i.e., less
232  influence of the differences in expected free energy across policies). For more details
233  regarding the associated mathematics, see (KJ Friston, Lin, et al., 2017; KJ Friston, Parr, et
234  al., 2017).
235

236       The 'A' matrix indicates which observations are generated by each

237  combination of hidden states (i.e., the likelihood mapping specifying the probability

238  that a particular set of observations would be observed given a particular set of

239  hidden states). The 'B' matrix is a transition matrix, indicating the probability that

240  one hidden state will evolve into another over time. The agent, based on the selected

241  policy, controls some of these transitions (e.g., those that pertain to the positions of

242  its limbs). The 'D' matrix encodes prior expectations about the initial hidden state

243  the agent will occupy. Finally, the 'C' matrix specifies prior preferences over

244  observations; it quantifies the degree to which different observed outcomes are

245  rewarding or punishing to the agent. In these models, observations and hidden

246  states can be factorized into multiple outcome *modalities* and hidden state *factors*.

247  This means that the likelihood mapping (the 'A' matrix) can also model the

248    interactions among different hidden states when generating outcomes

249    (observations).

250

251    **From principles to process theories**

252

253         One potential empirical advantage of the present approach stems from the

254    fact that active inference models have a plausible biological basis that affords

255    testable neurobiological predictions. Specifically, these models have well-articulated

256    companion micro-anatomical neural process theories, based on commonly used

257    message-passing algorithms (KJ Friston, FitzGerald, et al., 2017; Parr & Friston,

258    2018; Parr, Markovic, Kiebel, & Friston, 2019). In these process theories, for

259    example, the activation level of different neural populations (typically portrayed as

260    consisting of different cortical columns) can encode posterior probability estimates

261    over different hidden states. These activation levels can then be updated by synaptic

262    inputs with particular weights that convey the conditional probabilities encoded in

263    the 'A' and 'B' (among other) matrices described above, where active learning then

264    corresponds to associative synaptic plasticity. Phasic dopamine responses also play

265    a particular role in these models, by reporting changes in policy precision (i.e., the

266    degree of confidence in one policy over others) upon new observations (see Figure 2

267    and the associated legend for more details). In what follows, we describe how the

268    type of generative model – that underwrites these processes – was specified to

269    perform concept inference/learning.
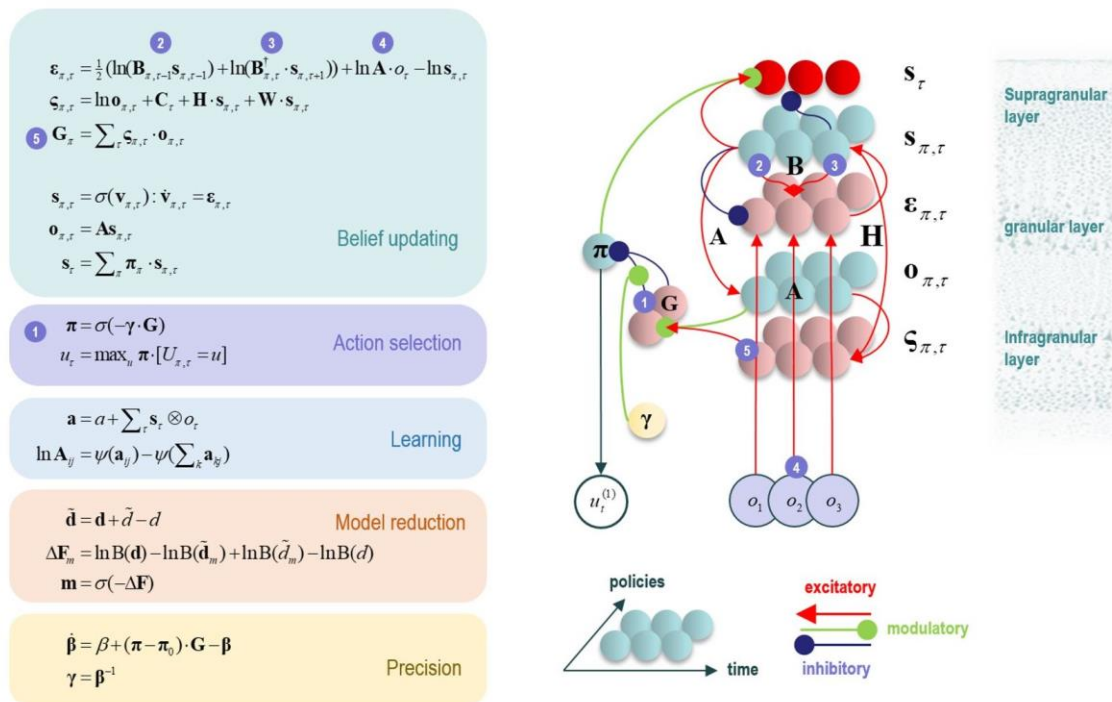
270

271



272

273 Figure 2. This figure illustrates the mathematical framework of active inference and
274 associated neural process theory used in the simulations described in this paper. The
275 differential equations in the left panel approximate Bayesian belief updating within the
276 graphical model depicted in the right panel of Figure 1 via a gradient descent on free energy
277 (**F**). The right panel also illustrates the proposed neural basis by which neurons making up
278 cortical columns could implement these equations. The equations have been expressed in
279 terms of two types of prediction errors. State prediction errors (**ε**) signal the difference
280 between the (logarithms of) expected states (**s**) under each policy and time point—and the
281 corresponding predictions based upon outcomes/observations (**A** matrix) and the
282 (preceding and subsequent) hidden states (**B** matrix, and, although not written, the **D**
283 matrix for the initial hidden states at the first time point). These represent prior and
284 likelihood terms respectively – also marked as messages 2, 3, and 4, which are depicted as
285 being passed between neural populations (colored balls) via particular synaptic
286 connections in the right panel. These (prediction error) signals drive depolarization (**v**) in
287 those neurons encoding hidden states (**s**), where the probability distribution over hidden
288 states is then obtained via a softmax (normalized exponential) function (**σ**). Outcome
289 prediction errors (**ς**) instead signal the difference between the (logarithms of) expected
290 observations (**o**) and those predicted under prior preferences (**C**). This term additionally
291 considers the expected ambiguity or conditional entropy (**H**) between states and outcomes
292 as well as a novelty term (**W**) reflecting the degree to which beliefs about how states
293 generate outcomes would change upon observing different possible state-outcome
294 mappings (computed from the **A** matrix). This prediction error is weighted by the expected
295 observations to evaluate the expected free energy (**G**) for each policy (**π**), conveyed via
296 message 5. These policy-specific free energies are then integrated to give the policy
297 expectations via a softmax function, conveyed through message 1. Actions at each time

298   point (**u**) are then chosen out of the possible actions under each policy (**U**) weighted by the
299   value (negative expected free energy) of each policy. In our simulations, the model learned
300   associations between hidden states and observations (**A**) via a process in which counts
301   were accumulated (**a**) reflecting the number of times the agent observed a particular
302   outcome when she believed that she occupied each possible hidden state. Although not
303   displayed explicitly, learning prior expectations over initial hidden states (**D**) is similarly
304   accomplished via accumulation of concentration parameters (**d**). These prior expectations
305   reflect counts of how many times the agent believes it previously occupied each possible
306   initial state. Concentration parameters are converted into expected log probabilities using
307   digamma functions ($\psi$). The way in which Bayesian model reduction was performed in this
308   paper is also written in the lower left (where B indicates a beta function, and **m** is the
309   posterior probability of each model). Here, the posterior distribution over initial states (**d**)
310   is used to assess the difference in the evidence ($\Delta$**F**) it provides for the number of hidden
311   states in the current model and other possible models characterized by fewer hidden states.
312   Prior concentration parameters are shown in italics, posterior in bold, and those priors and
313   posteriors associated with the reduced model are equipped with a tilde (**~**). As already
314   stated, the right panel illustrates a possible neural implementation of the update equations
315   in the left panel. In this implementation, probability estimates have been associated with
316   neuronal populations that are arranged to reproduce known intrinsic (within cortical area)
317   connections. Red connections are excitatory, blue connections are inhibitory, and green
318   connections are modulatory (i.e., involve a multiplication or weighting). These connections
319   mediate the message passing associated with the equations in the left panel. Cyan units
320   correspond to expectations about hidden states and (future) outcomes under each policy,
321   while red states indicate their Bayesian model averages (i.e., a "best guess" based on the
322   average of the probability estimates for the states and outcomes across policies, weighted
323   by the probability estimates for their associated policies. Pink units correspond to (state
324   and outcome) prediction errors that are averaged to evaluate expected free energy and
325   subsequent policy expectations (in the lower part of the network). This (neural) network
326   formulation of belief updating means that connection strengths correspond to the
327   parameters of the generative model described in the text. Learning then corresponds to
328   changes in the synaptic connection strengths. Only exemplar connections are shown to
329   avoid visual clutter. Furthermore, we have just shown neuronal populations encoding
330   hidden states under two policies over three time points (i.e., two transitions), whereas in
331   the task described in this paper there are greater number of allowable policies. For more
332   information regarding the mathematics and processes illustrated in this figure, see (KJ
333   Friston, Lin, et al., 2017; KJ Friston, Parr, et al., 2017).
334

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

336       To model concept inference, we constructed a simple task for an agent to

337   perform (see figure 1, left panel). In this task, the agent was presented with different

338   animals on different trials and asked to answer a question about the type of animal

339   that was seen. As described below, in some simulations the agent was asked to

340   report the type of animal that was learned previously; in other simulations, the

341    agent was instead asked a question that required conceptual generalization.

342    Crucially, to answer these questions the agent was required to observe different

343    animal features, where the identity of the animal depended on the combination of

344    features. There were three feature categories (size, color, and species-specific;

345    described further below) and two discrete time points in a trial (observe and

346    report).

347            To simulate concept learning (based on the task described above) we needed

348    to specify an appropriate generative model. Once this model has been specified, one

349    can use standard (variational) message passing to simulate belief updating and

350    behavior in a biologically plausible way: for details, please see (KJ Friston,

351    FitzGerald, et al., 2017; KJ Friston, Parr, et al., 2017). In our (minimal) model, the

352    first hidden state factor included (up to) eight levels, specifying four possible types

353    of birds and four possible types of fish (Figure 3A). The outcome modalities

354    included: a feature space including two size features (big, small), two color features

355    (gray, colorful), and two species-differentiating features (wings, gills). The 'A' matrix

356    specified a likelihood mapping between features and animal concepts, such that

357    each feature combination was predicted by an animal concept (Hawk, Pigeon,

358    Parrot, Parakeet, Sturgeon, Minnow, Whale shark, Clownfish). This model was

359    deliberately simple to allow for a clear illustration, but it is plausibly scalable to

360    include more concepts and a much larger feature space. The 'B' matrix for the first

361    hidden state factor was an identity matrix, reflecting the belief that the animal

362    identity was conserved during each trial (i.e., the animals were not switched out
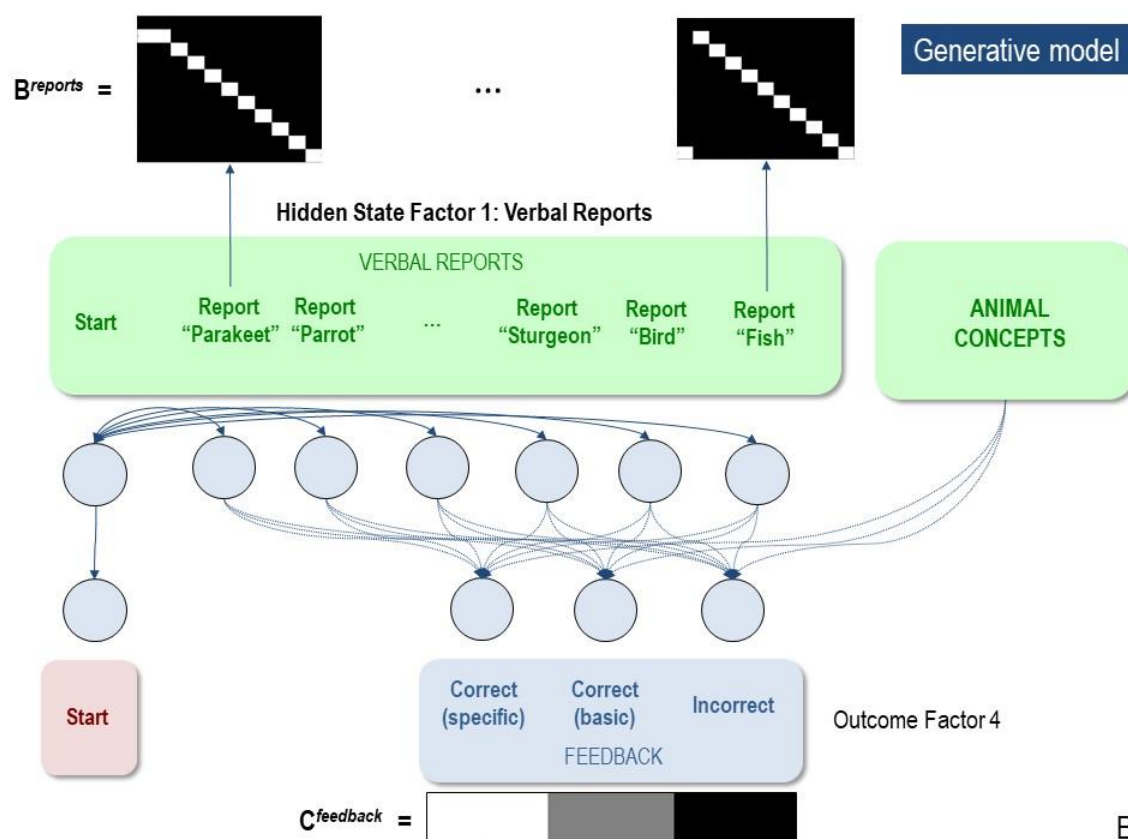
363    mid-trial).

Generative model

**Hidden State Factor 1: Concept Inference**

B Concepts

A Concepts

ANIMAL CONCEPTS

Parakeet   Parrot   Pigeon   Hawk   Clownfish   Whale Shark   Minnow   Sturgeon

Large   Small       Colorful   Gray        Wings   Gills

Size              Color              Species-specific

Observable Features

Outcome Factors 1, 2, 3

A

364
365

366

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

380

381 The second hidden state factor was the agent's report. That this is assumed

382 to factorise from the first hidden state factor means that there is no prior constraint

383 that links the chosen report to the animal generating observations. The agent could

384 report each of the eight possible specific animal categories, or opt for a less specific

385 report of a bird or a fish. Only one category could be reported at any time. Thus, the

386    agent had to choose to report only bird vs. fish or to report a more specific category.

387    In other words, the agent could decide upon the appropriate level of coarse-graining

388    of her responses (figure 3B).

389        During learning trials, the policy space was restricted such that the agent

390    could not provide verbal reports or observe corrective feedback (i.e., all it could do

391    is "stay still" in its initial state and observe the feature patterns presented). This

392    allowed the agent to learn concepts in an unsupervised manner (i.e. without being

393    told what the true state was or whether it was correct or incorrect). After learning,

394    active reporting was enabled, and the 'C' matrix was set so that the agent preferred

395    to report correct beliefs. We defined the preferences of the agent such that she

396    preferred correctly reporting specific category knowledge and was averse to

397    incorrect reports. This ensured that she only reported the general category of bird

398    vs. fish, unless sufficiently certain about the more specific category.

399        In the simulations reported below, there were two time points in each trial of

400    categorisation or conceptual inference. At the first time point, the agent was

401    presented with the animals features, and always began in a state of having made no

402    report (the "start" state). The agent's task was simply to observe the features, infer

403    the animal identity, and then report it (i.e., in reporting trials). Over 32 simulations

404    (i.e., 4 trials per animal), we confirmed that, if the agent already started out with full

405    knowledge of the animal concepts (i.e., a fully precise 'A' matrix), it would report the

406    specific category correctly 100% of the time. Over an additional 32 simulations, we

407    also confirmed that, if the agent was only equipped with knowledge of the

408    distinction between wings and gills (i.e., by replacing the rows in the 'A' matrix

409    corresponding to the mappings from animals to size and color with flat

410    distributions), it would report the generic category correctly 100% of the time but

411    would not report the specific categories.[2] This was an expected and straightforward

412    consequence of the generative model – but provides a useful example of how agents

413    trade off preferences and different types of uncertainty.

414

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

416

417         Having confirmed that our model could successfully recognize animals if

418    equipped with the relevant concepts (i.e., likelihood mappings) – we turn now to

419    concept learning.

420

421    **Concept acquisition**

422         We first examined our model's ability to acquire concept knowledge in two

423    distinct ways. This included 1) the agent's ability to "expand" (i.e., fill in an unused

424    column within) its state space and add new concepts, and 2) the agent's ability to

425    increase the granularity of its conceptual state space and learn more specific

426    concepts, when it already possessed broader concepts.

427

428    *Adding Concepts*

---

[2] 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 she disliked reporting incorrectly).

429     To assess whether our agent could expand her state space by acquiring a new

430     concept, we first set one column of the previously described model's 'A' matrix

431     (mapping an animal concept to its associated features) to be a uniform distribution[3];

432     creating an imprecise likelihood mapping for one concept – essentially, that concept

433     predicted all features with nearly equal probability. Here, we chose sturgeon (large,

434     gray, gills) as the concept for which the agent had no initial knowledge (see Figure

435     4A, right-most column of left-most 'pre-learning' matrix). We then generated 2000

436     observations based on the outcome statistics of a model with full knowledge of all

437     eight animals (the "generative process"), to test whether the model could learn the

438     correct likelihood mapping for sturgeon (note: this excessive number of

439     observations was used for consistency with later simulations, in which more

440     concepts had to be learned, and also to evaluate how performance improved as a

441     function of the number of observations the agent was exposed to; see figure 4B).

442     In these simulations, learning was implemented via updating (concentration)

443     parameters for the model's 'A' matrix after each trial. For details of these free energy

444     minimizing learning processes, please see (KJ Friston et al., 2016) as well as the left

445     panel of Figure 2 and associated legend. An intuitive way to think about this belief

446     updating process is that the strength of association between a concept and an

447     observation is quantified simply by counting how often they are inferred to co-

448     occur. This is exactly the same principle that underwrites Hebbian plasticity and

449     long-term potentiation (Brown, Zhao, & Leung, 2010). Crucially, policies were

---

[3] 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.

450    restricted during learning, such that the agent could not select reporting actions;

451    thus, learning was driven entirely by repeated exposure to different feature

452    combinations. We evaluated successful learning in two ways. First, we compared the

453    'A' matrix learned by the model to that of the generative process. Second, we

454    disabled learning after various trial numbers (i.e., such that concentration

455    parameters no longer accumulated) and enabled reporting. We then evaluated

456    reporting accuracy with 20 trials for each of the 8 concepts.

457         As shown in Figure 4A, the 'A' matrix (likelihood) mapping – learned by the

458    agent – and the column for sturgeon in particular, strongly resembled that of the

459    generative process. When first evaluating reporting, the model was 100 % accurate

460    across 20 reporting trials, when exposed to a sturgeon (reporting accuracy when

461    exposed to each of the other animals also remained at 100%) and first reached this

462    level of accuracy after around 50 exposures to all 8 animals (with equal probability)

463    (figure 4B). The agent also always chose to report specific categories (i.e., it never

464    chose to only report bird or fish). Model performance was stable over 8 repeated
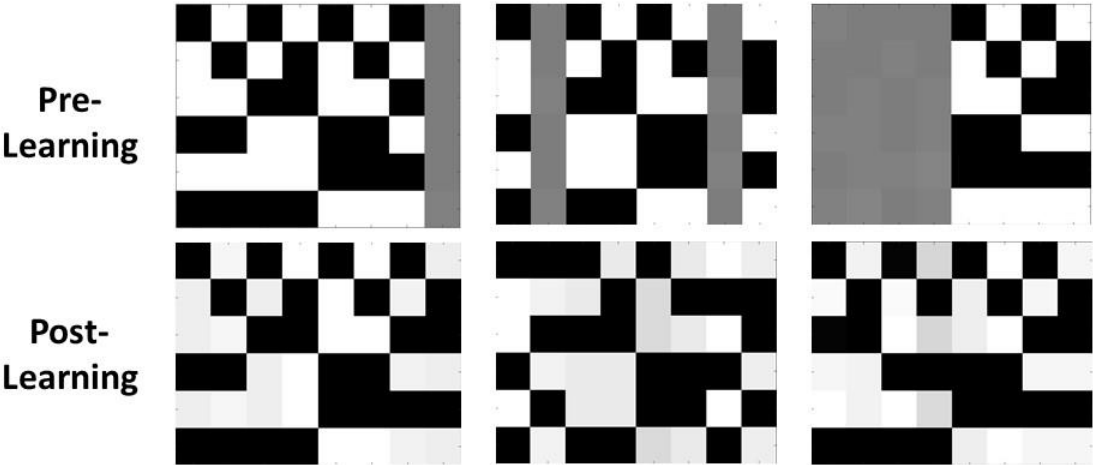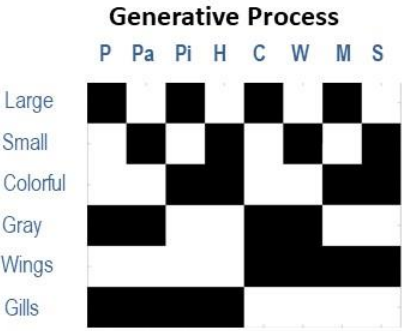
465    simulations.

466         Crucially, during learning, the agent was not told which state was generating

467    its observations. This meant that it had to solve both an inference and a learning

468    problem. First, it had to infer whether a given feature combination was better

469    explained by an existing concept, or by a concept that predicts features uniformly. In

470    other words, it had to decide that the features were sufficiently different – from

471    things it had seen before – to assign it a new hypothetical concept. Given that a novel

472    state is only inferred when another state is not a better explanation, this precludes

473    learning 'duplicate' states that generate the same patterns of observations. The

474    second problem is simpler. Having inferred that these outcomes are caused by

475    something new, the problem becomes one of learning a simple state-outcome

476    mapping through accumulation of Dirichlet parameters.

477        To examine whether this result generalized, we repeated these simulations

478    under conditions in which the agent had to learn more than one concept. When the

479    model needed to learn one bird (parakeet) and one fish (minnow), the model was

480    also able to learn the appropriate likelihood mapping for these 2 concepts (although

481    note that, because the agent did not receive feedback about the state it was in during

482    learning, the new feature mappings were often not assigned to the same columns as

483    in the generative process; see figure 4A). Reporting also reached 100% accuracy,

484    but required a notably greater number of trials. Across 8 repeated simulations, the

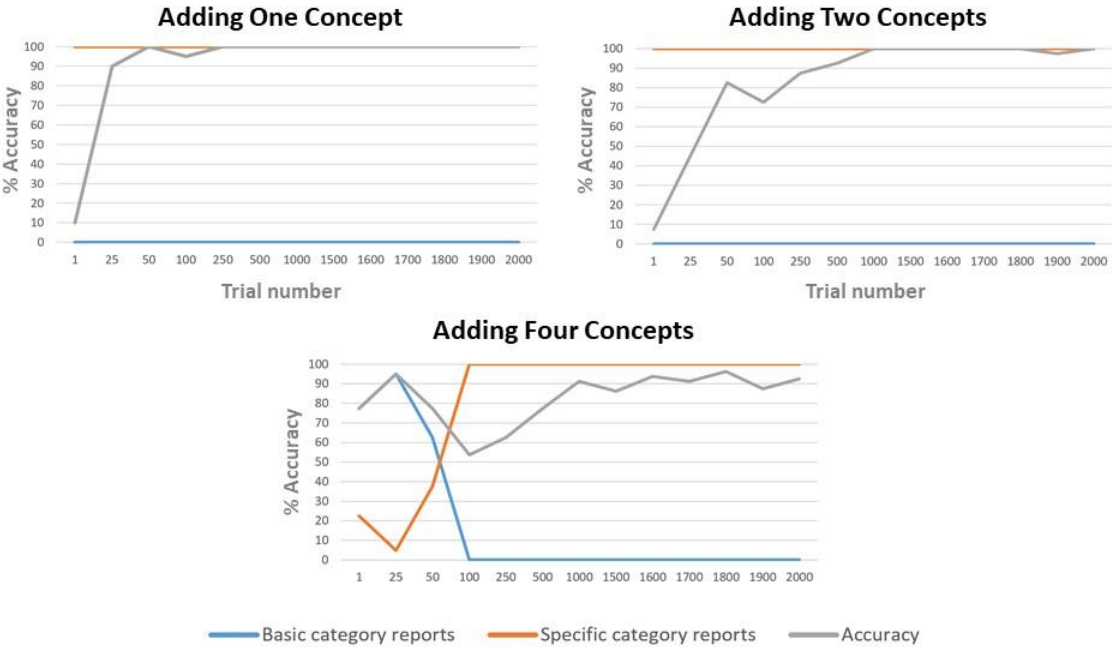485    mean accuracy reached by the model after 2000 trials was 98.75% (SD = 2%).

A



486

B



487

488     Figure 4. (A) illustration of representative simulation results in which the agent successfully
489     learned 1, 2, or 4 new animal concept categories with no prior knowledge beforehand. The
490     generative process is shown in the upper right, illustrating the feature combinations to be
491     learned. Pre-learning, either 1, 2 or 4 columns in the likelihood mapping began as a flat
492     distribution with a slight amount of Gaussian noise. The agent was then provided with 2000
493     observations of the 8 animals with equal probability. Crucially, the agent was prevented
494     from providing verbal reports during these 2000 trials and thus did not receive feedback
495     about the true identity of the animal. Thus, learning was driven completely by repeated
496     exposure in an unsupervised manner. Also note that, while the agent was successful at
497     learning the new concepts, it did not always assign the new feature patterns to the same
498     columns as illustrated in the generative process. This is to be expected given that the agent
499     received no feedback about the true hidden state that generated her observations. (B)
500     illustration of how reporting accuracy, and the proportion of basic category and specific
501     category responses, changed as a function of repeated exposures. This was accomplished by
502     taking the generative model at a number of representative trials and then testing it with 20
503     observations of each animal in which reporting was enabled. As can be seen, maximal
504     accuracy was achieved much more quickly when the agent had to learn fewer concepts.
505     When it had learned 4 concepts, it also began by reporting the general categories and then
506     subsequently became sufficiently confident to report the more specific categories.
507

508         When the model needed to learn all 4 birds, performance varied somewhat

509     more when the simulations were repeated. The learned likelihood mappings after

510     2000 trials always resembled that of the generative process, but with variable levels

511     of precision; notably, the model again assigned different concepts to different

512     columns relative to the generative process, as would be expected when the agent is

513     not given feedback about the state she is in. Over 8 repeated simulations, the model

514     performed well in 6 (92.50 % – 98.8 % accuracy) and failed to learn one concept in

515     the other 2 (72.50 % accuracy in each) due to overgeneralization (e.g., mistaking

516     parrot for Hawk in a majority of trials; i.e., the model simply learned that there are

517     large birds). Figure 4A and 4B illustrate representative results when the model was

518     successful (note: the agent never chose to report basic categories in the simulations

519     where only 1 or 2 concepts needed to be learned).

520        To further assess concept learning, we also tested the agent's ability to

521    successfully avoid state duplication. That is, we wished to confirm that the model

522    would only learn a new concept if actually presented with a new animal for which it

523    did not already have a concept. To do so, we equipped the model with knowledge of

524    seven out of the eight concept categories, and then repeatedly exposed it only to the

525    animals it already knew over 80 trials. We subsequently exposed it to the eighth

526    animal (Hawk) for which it did not already have knowledge over 20 additional

527    trials. As can be seen in figure 5, the unused concept column was not engaged during

528    the first 80 trials (bottom left and middle). However, in the final 20 trials, the agent

529    correctly inferred that her current conceptual repertoire was unable to explain her

530    new pattern of observations, leading the unused concept column to be adumbrated

531    and the appropriate state-observation mapping to be learned (bottom right). We

532    repeated these simulations under conditions in which the agent already had

533    knowledge of six, five, or four concepts. In all cases, we observed that unused

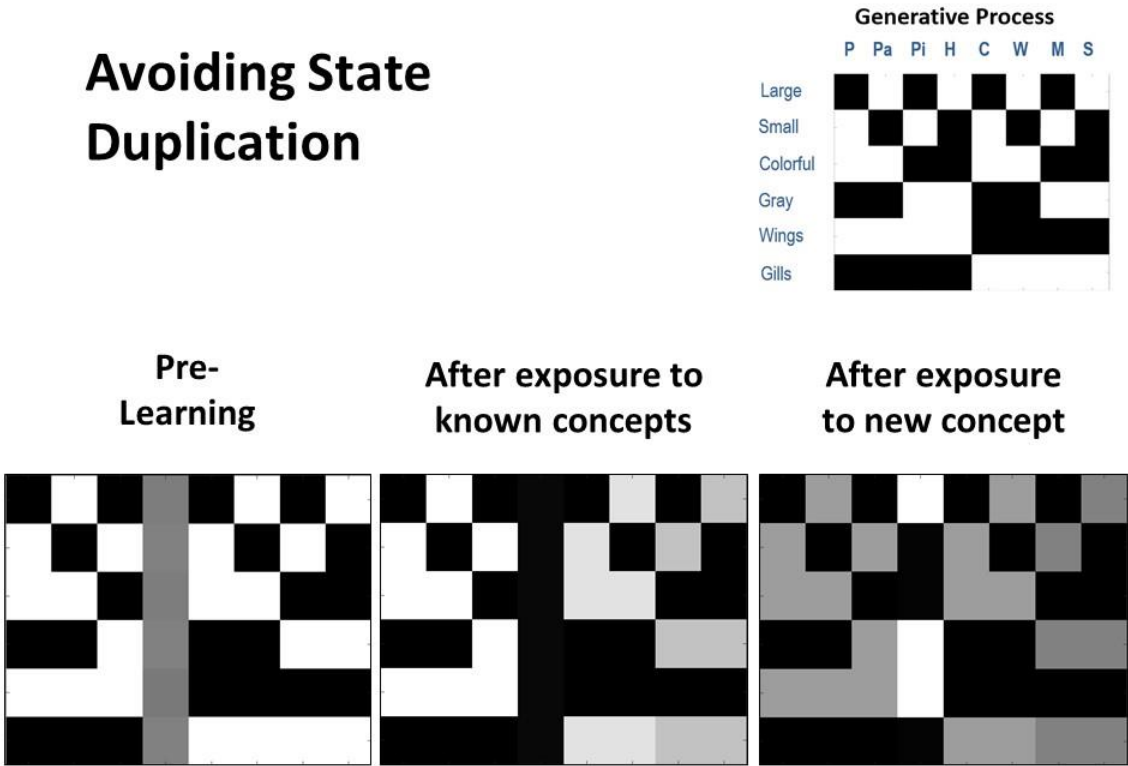534    concept columns were never engaged inappropriately.

535

536
537 Figure 5. Illustration of representative simulation results when the agent had to avoid
538 inappropriately learning a new concept (i.e., avoid state duplication) after only being
539 exposed to animals for which it already had knowledge. Here the agent began with prior
540 knowledge about seven concept categories and was also equipped with an eighth column
541 that could be engaged to learn a new concept category (bottom left). The agent was then
542 presented with several instances of each of the seven animals that she already knew (80
543 trials in total). In this simulation, the agent was successful in assigning each stimulus to an
544 animal concept she had already acquired and did not engage the unused concept 'slot'
545 (bottom middle). Finally, the agent was presented with a new animal (a hawk) that she did
546 not already know over 20 trials. In this case, the agent successfully engaged the additional
547 column (i.e., she inferred that none of the concepts she possessed could account for her new
548 observations) and learned the correct state-observation mapping (bottom right).
549

550          Crucially, these simulations suggest that adaptive concept learning needs to

551     be informed by existing knowledge about other concepts, such that a novel concept

552     should only be learned if observations cannot be explained with existing conceptual

553     knowledge. Here, this is achieved via the interplay of inference and learning, such

554     that agents initially have to infer whether to assign an observation to an existing

555    concept, and only if this is not possible is an 'open slot' employed to learn about a

556    novel concept.

557

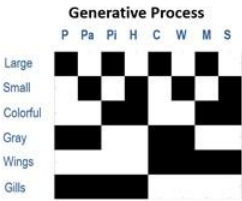558    *Increasing granularity*

559        Next, to explore the model's ability to increase the granularity of its concept

560    space, we first equipped the model with only the distinction between birds and fish

561    (i.e., the rows of the likelihood mapping corresponding to color and size features

562    were flattened in the same manner described above). We then performed the same

563    procedure used in our previous simulations. As can be seen in Figure 6A (bottom

564    left), the 'A' matrix learned by the model now more strongly resembled that of the

565    generative process. Figure 6B (bottom) also illustrates reporting accuracy and the

566    proportion of basic and specific category reports as a function of trial number. As

567    can be seen, the agent initially only reported general categories, and became

568    sufficiently confident to report specific categories after roughly 50 – 100 trials, but

569    her accuracy increased gradually over the next 1000 trials (i.e., the agent reported

570    specific categories considerably before its accuracy improved). Across 8 repeated

571    simulations, the final accuracy level reached was between 93% – 98% in 7

572    simulations, but the model failed to learn one concept in the 8th case, with 84.4%

573    overall accuracy (i.e., a failure to distinguish between pigeon and parakeet, and

574    therefore only learned a broader category of "small birds").

575        To assess whether learning basic categories first was helpful in subsequently

576    learning specific categories, we also repeated this simulation without any initial

577    knowledge of the basic categories. As exemplified in figure 6A and 6B, the model
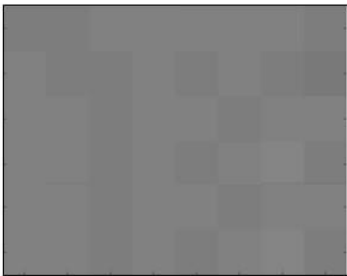
578     tended to perform reasonably well, but most often learned a less precise likelihood

579     mapping and reached a lower reporting accuracy percentage after 2000 learning

580     trials (across 8 repeated simulations: mean = 81.21%, SD = 6.39%, range from

581     68.80% – 91.30%). Thus, learning basic concept categories first appeared to

582     facilitate learning more specific concepts later.
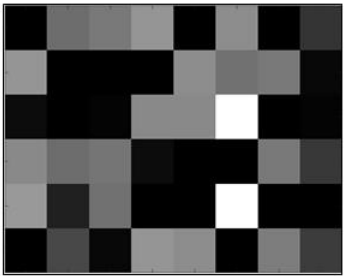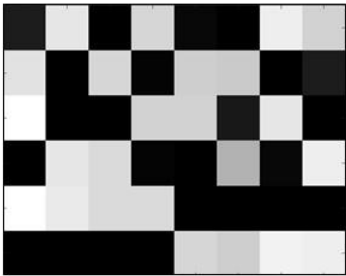


583
584

B

## Increasing Granularity:

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

**No Prior Knowledge**

**After Learning Basic Categories**

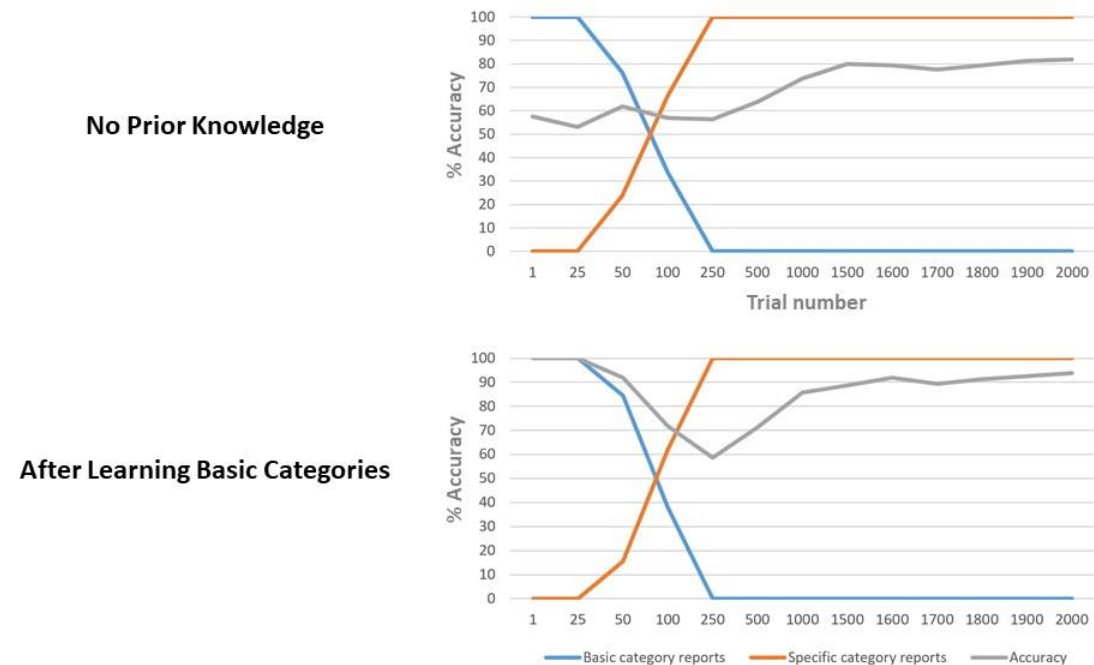Basic category reports — Specific category reports — Accuracy

Figure 6. (A, left) Illustration of representative simulation results when the agent had to learn to increase the granularity of her concept space. Here the agent began with prior knowledge about the basic concept categories (i.e., she 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 then 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 becomes 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, her final accuracy levels tended to be lower. In both cases, the agent began reporting specific categories before she achieved significant accuracy levels, therefore showing some initial overconfidence.

606       Overall, these findings provide a proof of principle that this sort of active

607    inference scheme can add concepts to a state space in an unsupervised manner (i.e.,

608    without feedback) based purely on (expected) free energy minimization. In this

609    case, it was able to accomplish this starting from a completely uninformative

610    likelihood distribution. It could also learn more granular concepts after already

611    acquiring more basic concepts, and our results suggest that learning granular

612    concepts may be facilitated by first learning basic concepts (e.g., as in currently

613    common educational practices).

614       The novel feature of this generative model involved 'building in' a number of

615    "reserve" hidden state levels. These initially had uninformative likelihood mappings;

616    yet, if a new pattern of features was repeatedly observed, and the model could not

617    account for this pattern with existing (informative) state-observation mappings,

618    these additional hidden state levels could be engaged to improve the model's

619    explanatory power. This approach therefore accommodates a simple form of

620    structure learning (i.e., model expansion).

621

622    **Integrating model expansion and reduction**

623

624    We next investigated ways in which model expansion could be combined with

625    Bayesian model reduction (KJ Friston, Lin, et al., 2017) – allowing the agent to adjust

626    her model to accommodate new patterns of observations, while also precluding

627    unnecessary conceptual complexity (i.e., over-fitting). To do so, we again allowed

628    the agent to learn from 2000 exposures to different animals as described in the

629     previous section – but also allowed the model to learn its 'D' matrix (i.e., accumulate

630     concentration parameters reflecting prior expectations over initial states). This

631     allowed an assessment of the agent's probabilistic beliefs about which hidden state

632     factor levels (animals) she had been exposed to. In different simulations, the agent

633     was only exposed to some animals and not others. We then examined whether a

634     subsequent model reduction step could recover the animal concepts presented

635     during the simulation; eliminating those concepts that were unnecessary to explain

636     the data at hand. The success of this 2-step procedure could then license the agent to

637     "reset" the unnecessary hidden state columns after concept acquisition, which

638     would have accrued unnecessary likelihood updates during learning. Doing so

639     would allow the optimal ability for those "reserve" states to be appropriately

640     engaged, if and when the agent was exposed to truly novel stimuli.

641          The 2nd step of this procedure was accomplished by applying Bayesian model

642     reduction to the 'D' matrix concentration parameters after learning. This is a form of

643     post-hoc model optimization (K. J. Friston et al., 2016; Karl Friston, Parr, & Zeidman,

644     2018) that rests upon estimation of a 'full' model, followed by analytic computation

645     of the evidence that would have been afforded to alternative models (with

646     alternative, 'reduced', priors) had they been used instead. Mathematically, this

647     procedure is a generalization of things like automatic relevance determination (Karl

648     Friston, Mattout, Trujillo-Barreto, Ashburner, & Penny, 2007; Wipf & Rao, 2007) or

649     the use of the Savage Dickie ratio in model comparison (Cornish & Littenberg,

650     2007). It is based upon straightforward probability theory and, importantly, has a

651     simple physiological interpretation; namely, synaptic decay and the elimination of
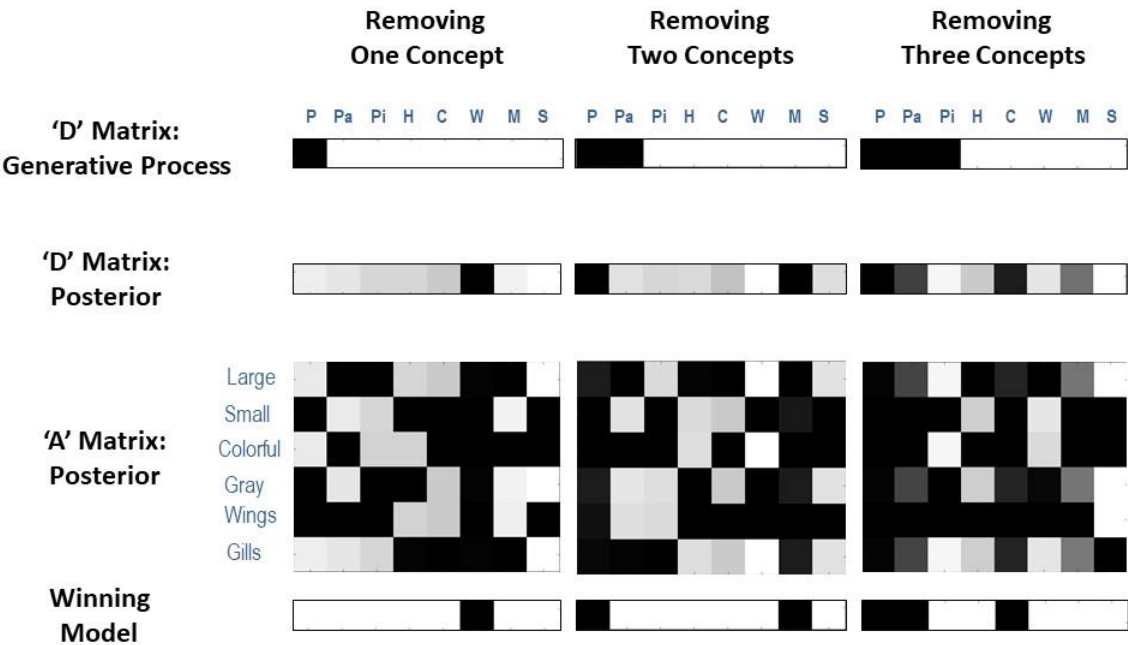
652    unused synaptic connections. In this (biological) setting, the concentration

653    parameters of the implicit Dirichlet distributions can be thought of as synaptic tags.

654    For a technical description of Bayesian model reduction techniques and their

655    proposed neural implementation, see (KJ Friston, Lin, et al., 2017; Hobson & Friston,

656    2012; Hobson, Hong, & Friston, 2014); see the left panel of Figure 2 for additional

657    details).

658         The posterior concentration parameters were compared to the prior

659    distribution for a full model (i.e., a flat distribution over 8 concepts) and prior

660    distributions for possible reduced models (i.e., which retained different possible

661    combinations of some but not all concepts; technically, reduced models were

662    defined such that the to-be-eliminated concepts were less likely than the to-be-

663    retained concepts). If Bayesian model reduction provided more evidence for one or

664    more reduced models, the reduced model with the most evidence was selected.

665    Note: an alternative would be to perform model reduction on the 'A' matrix, but this

666    is more complex due to a larger space of possible reduced models; it also does not

667    address the question of the number of hidden state levels to retain in a

668    straightforward manner.

669         In our first simulation, we presented our agent with all animals except for

670    parakeets with equal probability over 2000 trials. When compared to the full model,

671    the winning model corresponded to the correct 7-animal model matching the

672    generative process in 6/8 cases (log evidence differences ranged from -3.12 to -8.3),

673    and in 2/8 cases it instead selected a 6-animal model due to a failure to distinguish

674    between 2 specific concepts during learning (log evidence differences = -5.30, -

675    7.70). Figure 7 illustrates the results of a representative successful case). In the

676    successful cases, this would correctly license the removal of changes in the model's

677    'A' and 'D' matrix parameters for the 8th animal concept during learning in the

678    previous trials.  Similar results were obtained whenever any single animal type was

679    absent from the generative process.

## Bayesian Model Reduction After Learning



680

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

693

694       In a second simulation, the generative process contained 2 birds and all 4

695 fish. Here, the correct reduced model was correctly selected in 6/8 simulations (log

696 evidence differences range from -.96 to -8.24, with magnitudes greater than -3 in

697 5/6 cases), whereas it incorrectly selected the 5-animal model in 2 cases (log

698 evidence differences = -3.54, -4.50). In a third simulation, the generative process

699 contained 1 bird and all 4 fish. Here, the correct reduced model had the most

700 evidence in only 3/8 simulations (log evidence differences = -4.10, -4.11, -5.48),

701 whereas a 6-animal model was selected in 3/8 cases and a 3-animal and 7-animal

702 model were each selected once (log evidence differences > -3.0). Figure 7 also

703 illustrates representative examples of correct model recovery in these 2nd and 3rd

704 simulations.

705       While we have used the terms 'correct' and 'incorrect' above to describe the

706 model used to generate the data, we acknowledge that 'all models are wrong' (Box,

707 Hunter, & Hunter, 2005), and that the important question is not whether we can

708 recover the 'true' process used to generate the data, but whether we can arrive at

709 the simplest but accurate explanation for these data. The failures to recover the

710 'true' model highlighted above may reflect that a process other than that used to

711 generate the data could have been used to do so in a simpler way. Simpler here

712 means we would have to diverge to a lesser degree, from our prior beliefs, in order

713 to explain the data under a given model, relative to a more complex model. It is

714 worth highlighting the importance of the word *prior* in the previous sentence. This

715 means that the simplicity of the model is sensitive to our prior beliefs about it. To

716    illustrate this, we repeated the same model comparisons as above, but with precise

717    beliefs in an 'A' matrix that complies with that used to generate the data. Specifically,

718    we repeated the three simulations above but only enabled 'D' matrix learning (i.e.,

719    the model was already equipped with the 'A' matrix of the generative process). In

720    each case, Bayesian model reduction now uniquely identified the correct reduced

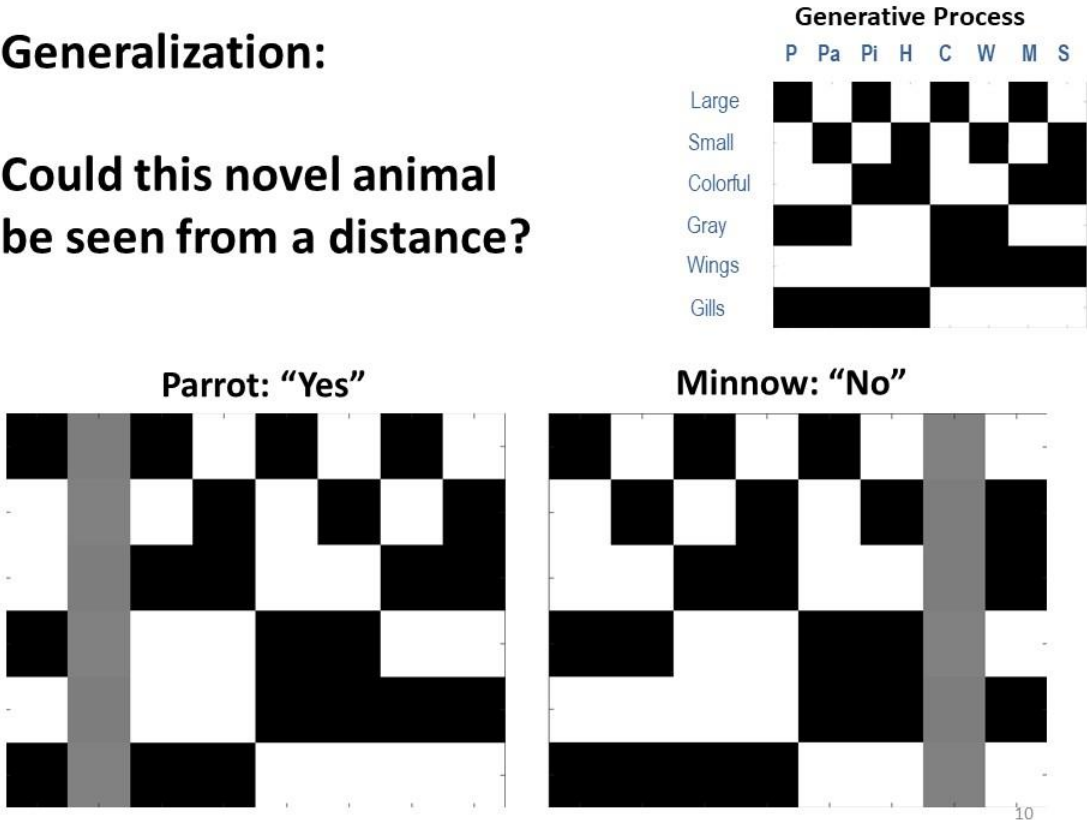721    model in 100% of cases.

722        These results demonstrate that – after a naïve model has expanded its hidden

723    state space to include likelihood mappings and initial state priors for a number of

724    concept categories – Bayesian model reduction can subsequently be used to

725    eliminate any parameter updates accrued for *one or two* redundant concept

726    categories. In practice, the 'A' and 'D' concentration parameters for these redundant

727    categories could be reset to their default pre-learning values – and could then be re-

728    engaged if new patterns of observations were repeatedly observed in the future.

729    However, when three concepts should have been removed, Bayesian model

730    reduction was much less reliable. This appeared to be due to imperfect 'A' matrix

731    learning, when occurring simultaneously with the (resultingly noisy) accumulation

732    of prior expectations over hidden states – as a fully precise 'A' matrix led to correct

733    model reduction in every case tested (i.e., suggesting that this type of model

734    reduction procedure could be improved by first allowing state-observation learning

735    to proceed alone, then subsequently allowing the model to learn prior expectations

736    over hidden states, which could then be used in model reduction).

737

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

739        One important ability afforded by concept learning is generalization. In a

740    final set of simulations, we asked if our model of concept knowledge could account

741    for generalization. To do so, we altered the model such that it no longer reported

742    what it saw, but instead had to answer a question that depended on generalization

743    from particular cross-category feature combinations. Specifically, the model was

744    shown particular animals and asked: "could this be seen from a distance?" The

745    answer to this question depended on both size and color, such that the answer was

746    yes only for colorful, large animals (i.e., assuming small or gray animals would blend

747    in with the sky or water and be missed).

748        Crucially, this question was asked of animals that the model had not been

749    exposed to, such that it had to generalize from knowledge it already possessed (see

750    Figure 8). To simulate and test for this ability, we equipped the model's 'A' matrix

751    with expert knowledge of 7 out of the 8 animals (i.e., as if these concepts had been

752    learned previously, as in our simulations above). The 8th animal was unknown to the

753    agent, in that it's likelihood mapping was set such that the 8th animal state "slot"

754    predicted all observations equally (i.e., with a small amount of Gaussian noise, as

755    above).  In one variant, the model possessed all concepts except for "parrot," and it

756    knew that the answer to the question was yes for "whale shark" but not for any

757    other concept it knew. To simulate one-shot generalization, learning was disabled

758    and a parrot (which it had never seen before) was presented 20 times to see if it

759    would correctly generalize and answer "yes" in a reliable manner. In another

760    variant, the model had learned all concepts except "minnow" and was tested the

761    same way to see if it would reliably provide the correct "no" response.

762    Here, we observed that in both of these cases (as well as all others we tested)

763    the model generalized remarkably well. It answered "yes" and "no" correctly in

764    100% of trials. Thus, the agent did not simply possess concepts to explain things it

765    saw. It instead demonstrated generalizable knowledge and could correctly answer

766    questions when seeing a novel stimulus.

767



768

769    Figure 8. Depiction of simulations in which we tested the agent's ability to generalize from
770    prior knowledge and correctly answered questions about new animals to which she had not
771    previously been exposed. In the simulations, the generative model was modified so that the
772    agent instead chose to report either "yes" or "no" to the question: "could this animal be seen
773    from a distance?" Here, the answer was only yes if the animal was both large and colorful.
774    We observed that when the agent started out with no knowledge of parrots it still correctly
775    answered this question 100% of the time, based only on its knowledge of other animals.
776    Similarly, when it started with no knowledge of minnows, it also correctly reported "no"
777    100% of the time. Thus, the agent was able to generalize from prior knowledge with no
778    additional learning.

779

780

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

782

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

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

802     observations can be sufficiently explained by existing conceptual knowledge – or

803     speak to the presence of a novel concept that motivates the use of an 'open slot'.

804     This process is cast as inference on (hidden) states. Second, if the agent infers that

805     there is a novel concept that explains current observations, it has to *learn* about the

806     specific feature configuration of that concept (i.e., novel state). This process

807     highlights the interplay between inference, which allows for the acquisition of

808     knowledge on a relatively short timescale, and learning, which allows for knowledge

809     acquisition on a longer and more stable timescale.

810             Similar considerations apply to the degree of 'similarity' of observations. In

811     the framework proposed here, we have assumed that the feature space of

812     observations is already learned and fixed. However, these feature spaces have to be

813     learned in the first place, which implies learning the underlying components or

814     feature dimensions that define observations. This relates closely to notions of

815     structure learning as dimensionality reduction based on covariance between

816     observations, as prominently discussed in the context of spatial navigation (Behrens

817     et al., 2018; Dordek, Soudry, Meir, & Derdikman, 2016; Stachenfeld et al., 2016;

818     Whittington, Muller, Mark, Barry, & Behrens, 2018).

819             Another important issue is how such abstract conceptual knowledge is

820     formed across different contexts or tasks. For example, the abstract concept of a

821     'bird' will be useful for learning about the fauna in a novel environment, but specific

822     types of birds – tied to a previous context – might be less useful in this regard. This

823     speaks to the formation of abstract, task-general knowledge that results from

824     training across different tasks, as recently discussed in the context of meta-

825    reinforcement learning (Ritter, Wang, Kurth-Nelson, & Botvinick, 2018; J X Wang et

826    al., 2016) with a putative link to the prefrontal cortex (Jane X. Wang et al., 2018). In

827    the present framework, such task-general knowledge would speak to the formation

828    of a hierarchical organisation that allows for the formation of conceptual knowledge

829    both within and across contexts. Also note that our proposed framework depends

830    on a pre-defined state space, including a pre-defined set of 'open slots' that allow for

831    novel context learning. The contribution of the present framework is to show how

832    these 'open slots' can be used for novel concept learning and be re-set based on

833    model reduction. It will be important to extend this approach towards learning the

834    structure of these models in the first place, including the appropriate number of

835    'open slots' (i.e., columns of the A-matrix) for learning in a particular content

836    domain and the relevant feature dimensions of observations (i.e., rows of A-matrix).

837    (Note: In addition to ontogenetic learning, in some cases structural priors regarding

838    the appropriate number of open slots [and relevant feature inputs for learning a

839    given state space of open slots] might also reflect inherited [i.e.,

840    genetically/developmentally pre-specified] patterns of structural neuronal

841    connectivity – based on what was adaptive within the evolutionary niche of a given

842    species – which could then be modified based on subsequent experience.)

843        This corresponds to a potentially powerful and simple application of

844    Bayesian model reduction, in which candidate models (i.e., reduced forms of a full

845    model) are readily identifiable based upon the similarity between the likelihoods

846    conditioned upon different hidden states. If two or more likelihoods are sufficiently

847    similar, the hidden states can be merged (by assigning the concentration

848    parameters accumulated during experience-dependent learning to one or other of

849    the hidden states). The ensuing change in model evidence scores the reduction in

850    complexity. If this reduction is greater than the loss of accuracy – in relation to

851    observations previously encountered – Bayesian model reduction will, effectively,

852    merge one state into another; thereby freeing up a state for the learning of new

853    concepts. We will demonstrate this form of structure learning via Bayesian model

854    reduction in future work.

855        We must also highlight here that cognitive science research on concept and

856    category learning has a rich empirical and theoretical history, including many

857    previously proposed formal models. While our primary focus has been on using

858    concept learning as an example of a more general approach by which state space

859    expansion and reduction can be implemented within future active inference

860    research, it is important to recognize this previous work and highlight where it

861    overlaps with the simulations we've presented. For example, our results suggesting

862    that first learning general categories facilitates the learning of more specific

863    categories relates to both classic and contemporary findings showing that children

864    more easily acquire "basic" and "superordinate" (e.g., dog, animal) concepts before

865    learning "subordinate" (e.g., chihuahua) concepts (Mervis & Rosch 1981; Murphy

866    2016), and that this may involve a type of "bootstrapping" process (Beck 2017).

867    Complementary work has also highlighted ways in which learning new words

868    during development can invoke a type of "placeholder" structure, which then

869    facilitates the acquisition of a novel concept – which bears some resemblance to our

870     notion of blank "concept slots" that can subsequently acquire meaningful semantics

871     (Gelman & Roberts 2017).

872          There is also a series of previously proposed formalisms within the literature

873     on category learning. For example, two previously proposed models – the "rational

874     model" (Anderson 1991; Sanborn et al. 2010) and the SUSTAIN model (Love et al.

875     2004) – both describe concept acquisition as involving cluster creation mechanisms

876     that depend on statistical regularities during learning and that use probabilistic

877     updating. The updating mechanisms within SUSTAIN are based on

878     surprise/prediction-error in the context of both supervised and unsupervised

879     learning. This model also down-weights previously created clusters when their

880     associated regularities cease to be observed in recent experience. Although not built

881     in intentionally, this type of mechanism emerges naturally within our model in two

882     ways. First, when a particular hidden state ceases to be inferred, concentration

883     parameters will accumulate to higher values for other hidden states in the **D** matrix,

884     reflecting relatively stronger prior expectations for hidden states that continue to be

885     inferred – which would favor future inference of those states over those absent from

886     recent experience. Second, if one pattern of observations were absent from recent

887     experience (while other patterns continued to be observed), concentration

888     parameters in the **A** matrix would also accumulate to higher values for patterns that

889     continued to be observed – resulting in relatively less confidence in the state-

890     outcome mapping for the less-observed pattern. (However, with respect to this

891     latter mechanism, so long as this mapping was sufficiently precise and distinct from

892     others [i.e., it had previously been observed many times farther in the past], this

893    would not be expected to prevent successful inference if this pattern were observed

894    again.)

895        It is also worth highlighting that, as our model is intended primarily as a

896    proof of concept and a demonstration of an available model expansion/reduction

897    approach that can be used within active inference research, it does not explicitly

898    incorporate some aspects – such as top-down attention – that are of clear

899    importance to cognitive learning processes, and that have been implemented in

900    previous models. For example, the adaptive resonance theory (ART) model

901    (Grossberg 1987) was designed to incorporate top-down attentional mechanisms

902    and feedback mechanisms to address a fundamental knowledge acquisition problem

903    – the temporal instability of previously learned information that can occur when a

904    system also remains sufficiently plastic to learn new (and potentially overlapping)

905    information. While our simulations do not explicitly incorporate these additional

906    complexities, there are clear analogues to the top-down and bottom-up feedback

907    exchange in ART within our model (e.g., the prediction and prediction-error

908    signaling within the neural process theory associated with active inference). ART

909    addresses the temporal instability problem primarily through mechanisms that

910    learn top-down expectancies that guide attention and match them with bottom-up

911    input patterns – which is quite similar to the prior expectations and likelihood

912    mappings used within active inference.

913        As an emergent property of the "first principles" approach in active

914    inference, our model therefore naturally incorporates the top-down effects in ART

915    simulations, which have been used to account for known context effects on

916    categorical perception within empirical studies (McClelland & Rumelhart 1981).

917    This is also consistent with more recent work on cross-categorization (Shafto et al.

918    2011), which has shown that human category learning is poorly accounted for by

919    both a purely bottom-up process (attempting to explain observed features) and a

920    purely top-down approach (involving attention-based feature selection) – and has

921    instead used simulations to show that a Bayesian joint inference model better fits

922    empirical data.

923        Other proposed Bayesian models of concept learning have also had

924    considerable success in predicting human generalization judgments (Goodman et al.

925    2008). The proof of concept model presented here has not been constructed to

926    explicitly compete with such models. It will be an important direction for future

927    work to explore the model's ability to scale up to handle more complex concept

928    learning problems. Here we simply highlight that the broadly Bayesian approach

929    within our model is shared with other models that have met with considerable

930    success – supporting the general plausibility of using this approach within active

931    inference research to model and predict the neural basis of these processes (see

932    below).

933

934                    **Potential advantages of the approach**

935        The present approach may offer some potential theoretical and empirical

936    advantages in comparison to previous work. One theoretical advantage corresponds

937    to the parsimony of casting this type of structure learning as an instance of Bayesian

938    model selection. When integrated with other aspects of the active inference

939     framework, this entails that perceptual inference, active learning, and structure

940     learning are all expressions of the same principle; namely, the minimization of

941     variational free energy, over three distinct timescales. A second, related theoretical

942     advantage is that, when this type of structure learning is cast as Bayesian model

943     selection/reduction, there is no need to invoke additional procedures or schemes

944     (e.g., nonparametric Bayes or 'stick breaking' processes; (S. Gershman & Blei,

945     2012)). Instead, a generative model with the capacity to represent a sufficiently

946     complex world will automatically learn causal structure in a way that contextualizes

947     active inference within active learning, and active learning within structure

948     learning.

949         Based on the process theories summarized in Figure 2, the present model

950     would predict that the brain contains "reserve" cortical columns and synapses (most

951     likely within secondary sensory and association cortices) available to capture new

952     patterns in observed features. To our knowledge, no direct evidence supporting the

953     presence of unused cortical columns in the brain has been observed, although the

954     generation of new neurons (with new synaptic connections) is known to occur in

955     the hippocampus (Chancey et al., 2013). "Silent synapses" have also been observed

956     in the brain, which does appear consistent with this prediction; such synapses can

957     persist into adulthood and only become activated when new learning becomes

958     necessary (e.g., see (Chancey et al., 2013; Funahashi, Maruyama, Yoshimura, &

959     Komatsu, 2013; Kerchner & Nicoll, 2008)). One way in which this idea of "spare

960     capacity" or "reserve" cortical columns might be tested in the context of

961     neuroimaging would be to examine whether greater levels of neural activation –

962    within conceptual processing regions – are observed after learning additional

963    concepts, which would imply that additional populations of neurons become

964    capable of being activated. In principle, single-cell recording methods might also test

965    for the presence of neurons that remain at baseline firing rates during task

966    conditions, but then become sensitive to new stimuli within the relevant conceptual

967    domain after learning.

968         Figure 9 provides a concrete example of two specific empirical predictions

969    that follow from simulating the neural responses that should be observed within our

970    concept learning task under these process theories. In the left panel, we plot the

971    firing rates (darker = higher firing rate) and local field potentials (rate of change in

972    firing rates) associated with neural populations encoding the probability of the

973    presence of different animals that would be expected across a number of learning

974    trials. In this particular example, the agent began with knowledge of the basic

975    categories of 'bird' and 'fish,' but needed to learn the eight more specific animal

976    categories over 50 interleaved exposures to each animal (only 10 equally spaced

977    learning trials involving the presentation of a parakeet are shown for simplicity). As

978    can be seen, early in learning the firing rates and local field potentials remain at

979    baseline levels; in contrast, as learning progresses, these neural responses take a

980    characteristic shape with more and more positive changes in firing rate in the

981    populations representing the most probable animal, while other populations drop

982    further and further below baseline firing rates.

983         The right panel depicts a similar simulation, but where the agent was

984    allowed to self-report what it saw on each trial (for clarity of illustration, we here

985    show 12 equally spaced learning trials for parakeet over 120 total trials). Enabling

986    policy selection allowed us to simulate expected phasic dopamine responses during

987    the task, corresponding to changes in the precision of the probability distribution

988    over policies after observing a stimulus on each trial. As can be seen, during early

989    trials the model predicts small firing rate increases when the agent is confident in its

990    ability to correctly report the more general animal category after observing a new

991    stimulus, and firing rate decreases when the agent becomes less confident in one

992    policy over others (i.e., as confidence in reporting the specific versus general

993    categories becomes more similar). Larger and larger phasic dopaminergic responses

994    are then expected as the agent becomes more and more confident in her ability to

995    correctly report the specific animal category upon observing a new stimulus. It will

996    be important for future neuroimaging studies to test these predictions in this type of

997    concept learning/stimulus categorization task.
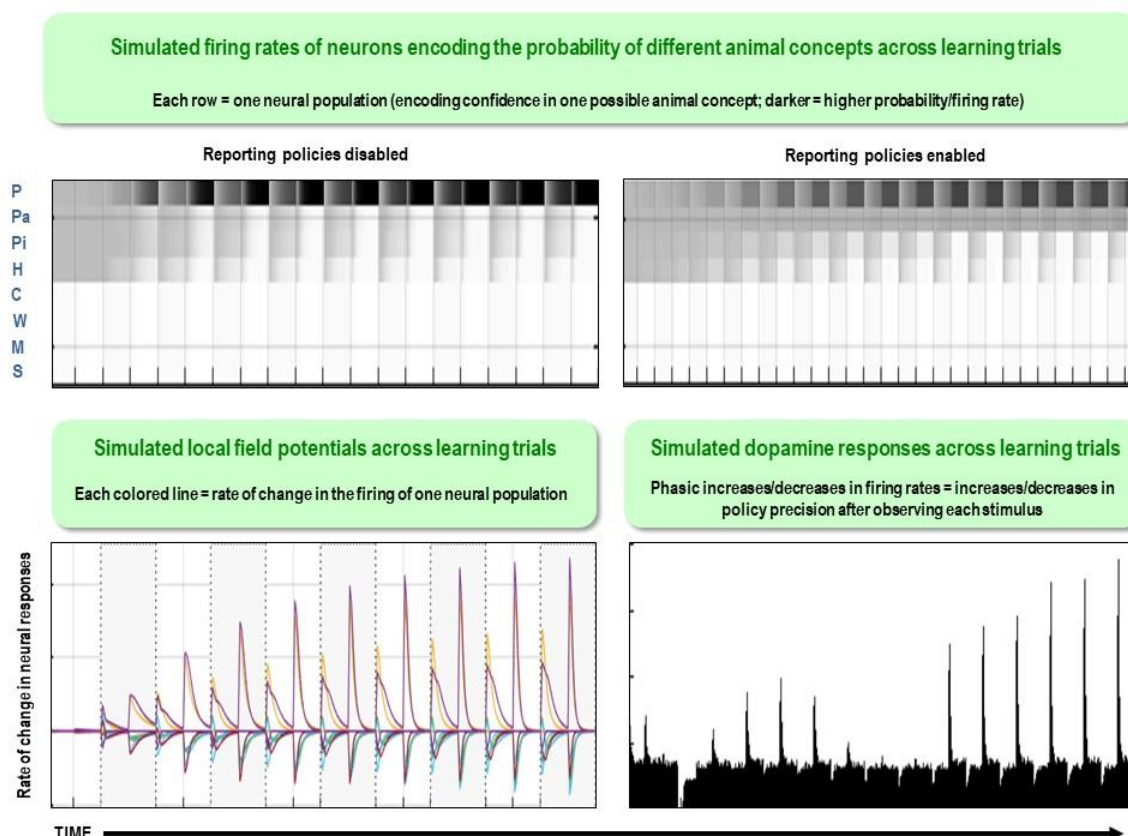
998

999

Figure 9. Simulated neuronal firing rates, local field potentials, and dopaminergic responses across learning trials based on the neural process theory associated with active inference that is summarized in Figure 2. The top left panel displays the predicted firing rates (darker = higher firing rate) of neural populations encoding the probability of each hidden state over 50 interleaved exposures to each animal (only 10 equally spaced learning trials involving the presentation of a parakeet are shown for simplicity) in the case where the agent starts out with knowledge of the basic animal categories but must learn the more specific categories. As can be seen, initially each of the four neural populations encoding possible bird categories (i.e., one row per possible category) have equally low firing rates (gray); as learning continues, firing rates increase for the 'parakeet' population and decrease for the others. The bottom left panel illustrates the predicted local field potentials (based on the rate of change in firing rates) that would be measured across the task. The top right panel displays the predicted firing rates of neural populations in an analogous simulation in which reporting policies were enabled (for clarity of illustration, we here show 12 equally spaced learning trials for parakeet over 120 total trials). Enabling policy selection allowed us to simulate the phasic dopaminergic responses (reporting changes in the precision of the probability distribution over policies) predicted to occur across learning trials; here the agent first becomes confident in her ability to correctly report the general animal category upon observing a stimulus, then becomes unsure about reporting specific versus general categories, and then becomes confident in her ability to report the specific categories.

1023

1024                                                            **Discussion**

1025

1026             The Active Inference formulation of concept learning presented here

1027     demonstrates a simple way in which a generative model can acquire both basic and

1028     highly granular knowledge of the hidden states/causes in its environment. In

1029     comparison to previous theoretical work using active inference (e.g., (M. Mirza,

1030     Adams, Mathys, & Friston, 2016; Parr & Friston, 2017; Schwartenbeck, FitzGerald,

1031     Mathys, Dolan, & Friston, 2015)), the novel aspect of our model was that it was

1032     further equipped with "reserve" hidden states initially devoid of content (i.e., these

1033     states started out with uninformative likelihood mappings that predicted all

1034     outcomes with roughly equal probability). Over multiple exposures to different

1035     stimuli, these hidden states came to acquire conceptual content that captured

1036     distinct statistical patterns in the features of those stimuli. This was accomplished

1037     via the model's ability to infer when its currently learned hidden states were unable

1038     to account for a new observation, leading an unused hidden state column to be

1039     engaged that could acquire a new state-observation mapping.

1040             Crucially, the model was able to start with some concepts and then expand its

1041     representational repertoire to learn others – but would only do so when a new

1042     stimulus was observed. This is conceptually similar to nonparametric Bayesian

1043     learning models, such as the "Chinese Room" process and the "Indian Buffet"

1044     process, that can also infer the need to invoke additional hidden causes with

1045     additional data (S. Gershman & Blei, 2012). These statistical learning models do not

1046    need to build in additional "category slots" for learning as in our model and can, in

1047    principle, entertain infinite state spaces. On the other hand, it is less clear at present

1048    how the brain could implement this type of learning. An advantage of our model is

1049    that learning depends solely on biologically plausible Hebbian mechanisms (for a

1050    possible neural implementation of model reduction, see (KJ Friston, Lin, et al., 2017;

1051    Hobson & Friston, 2012; Hobson et al., 2014)).

1052        The distinction between nonparametric Bayesian learning and the current

1053    active learning scheme may be important from a neurodevelopmental perspective

1054    as well. In brief, structure learning in this paper starts with a generative model with

1055    a type of structural prior reflecting a specific amount of built in 'spare capacity',

1056    where uncommitted or naive conceptual 'slots' are used to explain the sensorium,

1057    during optimization of free energy or model evidence. In contrast, nonparametric

1058    Bayesian approaches add new slots when appropriate. One might imagine that

1059    neonates are equipped with brains with 'spare capacity' (Baker & Tenenbaum,

1060    2014) that is progressively leveraged during neurodevelopment, much in the spirit

1061    of curriculum learning (Al-Muhaideb & Menai, 2011). This suggestion appears

1062    consistent with previous work demonstrating varying levels of category learning

1063    ability across the lifespan, which has previously been formally modeled as an

1064    individual difference in values of a parameter constraining the ability to form new

1065    clusters in response to surprising events (Love & Gureckis 2007) – which bears

1066    similarity to the idea of capacity limitations arising from finite numbers of concept

1067    slots in our model.

1068    In this sense, the current approach to structure learning may be better

1069    considered as active learning with generative models that are equipped with a large

1070    number of available hidden states capable of acquiring content, which are then

1071    judiciously reduced/reset – via a process of Bayesian model reduction.

1072    Furthermore, as in the acquisition of expertise, our model can also begin with broad

1073    category knowledge and then subsequently learn finer-grained within-category

1074    distinctions, which has received less attention from the perspective of the

1075    aforementioned models. Reporting broad versus specific category recognition is

1076    also a distinct aspect of our model – driven by differing levels of uncertainty and an

1077    expectation (preference) not to incorrectly report a more specific category.

1078    Our simulation results also demonstrated that, when combined with

1079    Bayesian model reduction, the model can guard against learning too many

1080    categories during model expansion – often retaining only the number of hidden

1081    causes actually present in its environment – and to keep "reserve" hidden states for

1082    learning about new causes if or when they appear. With perfect "expert" knowledge

1083    of the possible animal types it could observe (i.e., fully precise likelihood mappings

1084    matching the generative process) this was true in general. Interestingly, however,

1085    with an imperfectly learned likelihood mapping, model reduction only succeeded

1086    when the agent had to remove either 1 or 2 concepts from her model; when 3

1087    potential categories needed to be removed, the correct reduced model was

1088    identified less than half the time. It would be interesting to empirically test whether

1089    similar learning difficulties are present in humans.

1090    Neurobiological theories associated with Active Inference also make

1091    predictions about the neural basis of this process (Hobson & Friston, 2012; Hobson

1092    et al., 2014). Specifically, during periods of rest (e.g., daydreaming) or sleep, it is

1093    suggested that, because sensory information is down-weighted, learning is driven

1094    mainly by internal model simulations (e.g., as appears to happen in the phenomenon

1095    of hippocampal replay; (Feld & Born, 2017; Lewis, Knoblich, & Poe, 2018; Pfeiffer &

1096    Foster, 2013)); this type of learning can accomplish a model reduction process in

1097    which redundant model parameters are identified and removed to prevent model

1098    over-fitting and promote selection of the most parsimonious model that can

1099    successfully account for previous observations.  This is consistent with work

1100    suggesting that, during sleep, many (but not all) synaptic strength increases

1101    acquired in the previous day are attenuated (Tononi & Cirelli, 2014). The role of

1102    sleep and daydreaming in keeping "reserve" representational resources available

1103    for model expansion could therefore be especially important to concept learning –

1104    consistent with the known role of sleep in learning and memory (Ackermann &

1105    Rasch, 2014; Feld & Born, 2017; Perogamvros & Schwartz, 2012; Stickgold, Hobson,

1106    Fosse, & Fosse, 2001; Walker & Stickgold, 2010).

1107    In addition, an emergent feature of our model was its ability to generalize

1108    prior knowledge to new stimuli to which it had not previously been exposed. In fact,

1109    the model could correctly generalize upon a single exposure to a new stimulus – a

1110    type of "one-shot learning" capacity qualitatively similar to that observed in humans

1111    (Landau, Smith, & Jones, 1988; E. Markman, 1989; Xu & Tenenbaum, 2007b). While

1112    it should be kept in mind that the example we have provided is very simple, it

1113    demonstrates the potential usefulness of this novel approach. Some other

1114    prominent approaches in machine-learning (e.g., deep learning) tend to require

1115    larger amounts of data (Geman et al., 1992; Hinton et al., 2012; LeCun et al., 2015;

1116    Lecun et al., 1998; Mnih et al., 2015), and do not learn the rich structure that allows

1117    humans to use concept knowledge in a wide variety of generalizable functions

1118    (Barsalou, 1983; Biederman, 1987; Feldman, 1997; Jern & Kemp, 2013; A. B.

1119    Markman & Makin, 1998; Osherson & Smith, 1981; Ward, 1994; Williams &

1120    Lombrozo, 2010). Other recent hierarchical Bayesian approaches in cognitive

1121    science have made progress in this domain, however, by modeling concepts as types

1122    of probabilistic programs (Ghahramani, 2015; Goodman, Tenenbaum, &

1123    Gerstenberg, 2015; Lake et al., 2015).

1124            It is important to note that this model is deliberately simple and is meant

1125    only to represent a proof of principle that categorical inference and conceptual

1126    knowledge acquisition can be modeled within this particular neurocomputational

1127    framework, and to present this approach as a potentially useful tool in future active

1128    inference research. We chose a particular set of feature combinations to illustrate

1129    this, but it remains to be demonstrated that learning in this model would be equally

1130    successful with a larger feature space and set of learnable hidden causes. Due to

1131    limited scope, we have also not thoroughly addressed all other overlapping lines of

1132    research. For example, work on exemplar models of concepts has also led to other

1133    computational approaches. As one example, the EBRW model (Nosofsky & Palmeri

1134    1997) has demonstrated ways of linking exemplar learning to drift diffusion models.

1135    Another model within this line of research is the ALCOVE model (Nosofsky et al.

1136    1994) – an error-driven connectionist model of exemplar-based category learning

1137    that employs selective attention and learns attentional weights (this model also

1138    built on earlier work; see (Nosofsky 2011)). Yet another connectionist model with

1139    some conceptual overlap to our own is the DIVA model, which learns categories by

1140    recoding observations as task-constrained principle components and uses model fit

1141    for subsequent recognition (Kurtz 2007). It will be important in future work to

1142    examine the strengths and limitations of a scaled-up version of our approach in

1143    relation to these other models.

1144          Finally, another topic for future work would be the expansion of this type of

1145    model to context-specific learning (e.g., with an additional hidden state factor for

1146    encoding distinct contexts). In such cases, regularities in co-occurring features differ

1147    in different contexts and other cues to context may not be directly observable (e.g.,

1148    the same species of bird could be a slightly different color or size in different parts

1149    of the world that otherwise appear similar) – creating difficulties in inferring when

1150    to update previously learned associations and when to instead acquire competing

1151    associations assigned to new contexts. At present, it is not clear whether the

1152    approach we have illustrated would be successful at performing this additional

1153    function, although the process of inferring the presence of a new hidden state level

1154    in a second hidden state factor encoding context would be similar to what we have

1155    shown within a single state factor (for related work on context-dependent

1156    contingency learning, see (S. J. Gershman et al., 2017; S. Gershman, Jones, Norman,

1157    Monfils, & Niv, 2013)). Another point worth highlighting is that we have made

1158    particular choices with regard to various model parameters and the number of

1159    observations provided during learning. Further investigations of the space of these

1160    possible parameter settings will be important. With this in mind, however, our

1161    current modelling results could offer additional benefits. For example, the model's

1162    simplicity could be amenable to empirical studies of saccadic eye movements

1163    toward specific features during novel category learning (e.g. following the approach

1164    of (M. B. Mirza, Adams, Mathys, & Friston, 2018)). This approach could also be

1165    combined with measures of neural activity in humans or other animals, allowing

1166    more direct tests of the neural predictions highlighted above. In addition, the

1167    introduction of exploratory, novelty-seeking, actions could be used to reduce the

1168    number of samples required for learning, with agents selecting those data that are

1169    most relevant.

1170        In conclusion, the Active Inference scheme we have described illustrates

1171    feature integration in the service of conceptual inference: it can successfully

1172    simulate simple forms of concept acquisition and concept differentiation (i.e.

1173    increasing granularity), and it spontaneously affords one-shot generalization.

1174    Finally, it speaks to empirical work in which behavioral tasks could be designed to

1175    fit such models, which would allow investigation of individual differences in concept

1176    learning and its neural basis. For example, such a model can simulate (neuronal)

1177    belief updating to predict neuroimaging responses as we illustrated above; i.e., to

1178    identify the neural networks engaged in evidence accumulation and learning

1179    (Schwartenbeck et al., 2015). In principle, the model parameters (e.g., 'A' matrix

1180    precision) can also be fit to behavioral choices and reaction times – and thereby

1181    phenotype subjects in terms of the priors under which they infer and learn

1182  (Schwartenbeck & Friston, 2016). This approach could therefore advance

1183  neurocomputational approaches to concept learning in several directions.

1184

1185  **Software note**

1186  Although the generative model – specified by the various matrices described in this

1187  paper – changes from application to application, the belief updates are generic and

1188  can be implemented using standard routines (here **spm_MDP_VB_X.m**). These

1189  routines are available as Matlab code in the DEM toolbox of the most recent version

1190  of SPM academic software: http://www.fil.ion.ucl.ac.uk/spm/. The simulations in

1191  this paper can be reproduced (and customized) via running the Matlab code

1192  included here is supplementary material (**Concepts_model.m**).

1193

1194

1195  **References**

1196

1197  Ackermann, S., & Rasch, B. (2014). Differential Effects of Non-REM and REM Sleep

1198      on Memory Consolidation? *Current Neurology and Neuroscience Reports*, *14*(2),

1199      430. https://doi.org/10.1007/s11910-013-0430-8

1200  Al-Muhaideb, S., & Menai, M. E. B. (2011). Evolutionary computation approaches to

1201      the Curriculum Sequencing problem. *Natural Computing*, *10*(2), 891–920.

1202      https://doi.org/10.1007/s11047-010-9246-5

1203  Anderson, J. R. (1991). The Adaptive Nature of Human Categorization. *Psychological

1204      Review*, *98*(3), 409–429. https://doi.org/10.1037/0033-295X.98.3.409

1205  Baker, C., & Tenenbaum, J. (2014). Modeling human plan recognition using Bayesian

1206      theory of mind. In G. Sukthankar, C. Geib, H. Dui, D. Pynadath, & R. Goldman

1207      (Eds.), *Plan, activity, and intent recognition* (pp. 177–204). Boston: Morgan

1208      Kaufmann.

1209  Barsalou, L. W. (1983). Ad hoc categories. *Memory & Cognition*, *11*(3), 211–227.

1210      https://doi.org/10.3758/bf03196968

1211  Behrens, T. E. J., Muller, T. H., Whittington, J. C. R., Mark, S., Baram, A. B., Stachenfeld,

1212      K. L., & Kurth-Nelson, Z. (2018). What Is a Cognitive Map? Organizing

1213      Knowledge for Flexible Behavior. *Neuron*, *100*(2), 490–509.

1214      https://doi.org/10.1016/J.NEURON.2018.10.002

1215  Biederman, I. (1987). Recognition-by-components: a theory of human image

1216      understanding. *Psychological Review*, *94*(2), 115–147.

1217      https://doi.org/10.1037/0033-295X.94.2.115

1218  Botvinick, M. M., Niv, Y., & Barto, A. C. (2009). Hierarchically organized behavior and

1219      its neural foundations: A reinforcement learning perspective. *Cognition*, *113*(3),

1220      262–280. https://doi.org/10.1016/J.COGNITION.2008.08.011

1221  Box, G. E., Hunter, J. S., & Hunter, W. G. (2005). *Statistics for experimenters. Wiley*

1222      *Series in Probability and Statistics.* Hoboken, NJ.: Wiley.

1223  Brown, T. H., Zhao, Y., & Leung, V. (2010). Hebbian plasticity. In *Encyclopedia of*

1224      *Neuroscience* (pp. 1049–1056). https://doi.org/10.1016/B978-008045046-

1225      9.00796-8

1226  Chancey, J., Adlaf, E., Sapp, M., Pugh, P., Wadiche, J., & Overstreet-Wadiche, L. (2013).

1227      GABA depolarization is required for experience-dependent synapse unsilencing

1228    in adult-born neurons. *The Journal of Neuroscience : The Official Journal of the*

1229    *Society for Neuroscience*, *33*(15), 6614–6622.

1230    https://doi.org/10.1523/JNEUROSCI.0781-13.2013

1231  Conant, C., & Ashbey, W. (1970). Every good regulator of a system must be a model

1232    of that system. *International Journal of Systems Science*, *1*(2), 89–97.

1233    https://doi.org/10.1080/00207727008920220

1234  Cornish, N. J., & Littenberg, T. B. (2007). Tests of Bayesian model selection

1235    techniques for gravitational wave astronomy. *Physical Review D*, *76*(8), 083006.

1236    https://doi.org/10.1103/PhysRevD.76.083006

1237  Dordek, Y., Soudry, D., Meir, R., & Derdikman, D. (2016). Extracting grid cell

1238    characteristics from place cell inputs using non-negative principal component

1239    analysis. *ELife*, *5*(MARCH2016), 1–36. https://doi.org/10.7554/eLife.10094

1240  Feld, G., & Born, J. (2017). Sculpting memory during sleep: concurrent consolidation

1241    and forgetting. *Current Opinion in Neurobiology*, *44*, 20–27.

1242    https://doi.org/10.1016/J.CONB.2017.02.012

1243  Feldman, J. (1997). The Structure of Perceptual Categories. *Journal of Mathematical*

1244    *Psychology*, *41*(2), 145–170. https://doi.org/10.1006/jmps.1997.1154

1245  Friston, K. J., Litvak, V., Oswal, A., Razi, A., Stephan, K. E., van Wijk, B. C. M., …

1246    Zeidman, P. (2016). Bayesian model reduction and empirical Bayes for group

1247    (DCM) studies. *NeuroImage*, *128*, 413–431.

1248    https://doi.org/10.1016/J.NEUROIMAGE.2015.11.015

1249  Friston, Karl, Mattout, J., Trujillo-Barreto, N., Ashburner, J., & Penny, W. (2007).

1250    Variational free energy and the Laplace approximation. *NeuroImage*, *34*(1),

1251      220–234. https://doi.org/10.1016/J.NEUROIMAGE.2006.08.035

1252  Friston, Karl, Parr, T., & Zeidman, P. (2018). Bayesian model reduction. Retrieved

1253      from http://arxiv.org/abs/1805.07092

1254  Friston, Karl, & Penny, W. (2011). Post hoc Bayesian model selection. *NeuroImage*,

1255      *56*(4), 2089–2099. https://doi.org/10.1016/J.NEUROIMAGE.2011.03.062

1256  Friston, KJ. (2010). The free-energy principle: a unified brain theory? *Nature*

1257      *Reviews. Neuroscience*, *11*(2), 127–138. https://doi.org/10.1038/nrn2787

1258  Friston, KJ, FitzGerald, T., Rigoli, F., Schwartenbeck, P., O Doherty, J., & Pezzulo, G.

1259      (2016). Active inference and learning. *Neuroscience and Biobehavioral Reviews*,

1260      *68*, 862–879. https://doi.org/10.1016/j.neubiorev.2016.06.022

1261  Friston, KJ, FitzGerald, T., Rigoli, F., Schwartenbeck, P., & Pezzulo, G. (2017). Active

1262      Inference: A Process Theory. *Neural Computation*, *29*(1), 1–49.

1263      https://doi.org/10.1162/NECO_a_00912

1264  Friston, KJ, Lin, M., Frith, C., Pezzulo, G., Hobson, J., & Ondobaka, S. (2017). Active

1265      Inference, Curiosity and Insight. *Neural Computation*, *29*(10), 2633–2683.

1266      https://doi.org/10.1162/neco_a_00999

1267  Friston, KJ, Parr, T., & de Vries, B. (2017). The graphical brain: Belief propagation

1268      and active inference. *Network Neuroscience*, *1*(4), 381–414.

1269      https://doi.org/10.1162/NETN_a_00018

1270  Funahashi, R., Maruyama, T., Yoshimura, Y., & Komatsu, Y. (2013). Silent synapses

1271      persist into adulthood in layer 2/3 pyramidal neurons of visual cortex in dark-

1272      reared mice. *Journal of Neurophysiology*, *109*(8), 2064–2076.

1273      https://doi.org/10.1152/jn.00912.2012

1274    Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural Networks and the

1275        Bias/Variance Dilemma. *Neural Computation*, *4*(1), 1–58.

1276        https://doi.org/10.1162/neco.1992.4.1.1

1277    Gershman, S., & Blei, D. (2012). A tutorial on Bayesian nonparametric models.

1278        *Journal of Mathematical Psychology*, *56*(1), 1–12.

1279        https://doi.org/10.1016/J.JMP.2011.08.004

1280    Gershman, S. J., Monfils, M.-H., Norman, K. A., & Niv, Y. (2017). The computational

1281        nature of memory modification. *ELife*, *6*. https://doi.org/10.7554/eLife.23763

1282    Gershman, S. J., & Niv, Y. (2010). Learning latent structure: carving nature at its

1283        joints. *Current Opinion in Neurobiology*, *20*(2), 251–256.

1284        https://doi.org/10.1016/J.CONB.2010.02.008

1285    Gershman, S., Jones, C., Norman, K., Monfils, M., & Niv, Y. (2013). Gradual extinction

1286        prevents the return of fear: implications for the discovery of state. *Frontiers in*

1287        *Behavioral Neuroscience*, *7*, 164. https://doi.org/10.3389/fnbeh.2013.00164

1288    Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence.

1289        *Nature*, *521*(7553), 452–459. https://doi.org/10.1038/nature14541

1290    Goodman, N. D., Tenenbaum, J. B., & Gerstenberg, T. (2015). *Concepts: New*

1291        *Directions*. (E. Margolis & S. Laurence, Eds.). Cambridge, MA: MIT Press.

1292    Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A., Jaitly, N., … Kingsbury, B. (2012).

1293        Deep Neural Networks for Acoustic Modeling in Speech Recognition: The

1294        Shared Views of Four Research Groups. *IEEE Signal Processing Magazine*, *29*(6),

1295        82–97. https://doi.org/10.1109/MSP.2012.2205597

1296    Hobson, J., & Friston, K. (2012). Waking and dreaming consciousness:

1297    neurobiological and functional considerations. *Progress in Neurobiology*, *98*(1),

1298    82–98. https://doi.org/10.1016/j.pneurobio.2012.05.003

1299    Hobson, J., Hong, C.-H., & Friston, K. (2014). Virtual reality and consciousness

1300    inference in dreaming. *Frontiers in Psychology*, *5*, 1133.

1301    https://doi.org/10.3389/fpsyg.2014.01133

1302    Jern, A., & Kemp, C. (2013). A probabilistic account of exemplar and category

1303    generation. *Cognitive Psychology*, *66*(1), 85–125.

1304    https://doi.org/10.1016/j.cogpsych.2012.09.003

1305    Kemp, C., Perfors, A., & Tenenbaum, J. (2007). Learning overhypotheses with

1306    hierarchical Bayesian models. *Developmental Science*, *10*(3), 307–321.

1307    https://doi.org/10.1111/j.1467-7687.2007.00585.x

1308    Kerchner, G., & Nicoll, R. (2008). Silent synapses and the emergence of a

1309    postsynaptic mechanism for LTP. *Nature Reviews. Neuroscience*, *9*(11), 813–

1310    825. https://doi.org/10.1038/nrn2501

1311    Lake, B. B. M., Salakhutdinov, R., & Tenenbaum, J. J. B. (2015). Human-level concept

1312    learning through probabilistic program induction. *Science (New York, N.Y.)*,

1313    *350*(6266), 1332–1338. https://doi.org/10.1126/science.aab3050

1314    Landau, B., Smith, L. B., & Jones, S. S. (1988). The importance of shape in early lexical

1315    learning. *Cognitive Development*, *3*(3), 299–321.

1316    https://doi.org/10.1016/0885-2014(88)90014-7

1317    LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, *521*(7553), 436–

1318    444. https://doi.org/10.1038/nature14539

1319    Lecun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning

1320    applied to document recognition. *Proceedings of the IEEE*, *86*(11), 2278–2324.

1321    https://doi.org/10.1109/5.726791

1322  Lewis, P., Knoblich, G., & Poe, G. (2018). How Memory Replay in Sleep Boosts

1323    Creative Problem-Solving. *Trends in Cognitive Sciences*, *22*(6), 491–503.

1324    https://doi.org/10.1016/j.tics.2018.03.009

1325  Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). SUSTAIN: a network model of

1326    category learning. *Psychological Review*, *111*(2), 309–332.

1327    https://doi.org/10.1037/0033-295X.111.2.309

1328  MacKay, D. J. C., & Peto, L. C. B. (1995). A hierarchical Dirichlet language model.

1329    *Natural Language Engineering*, *1*(03), 289–308.

1330    https://doi.org/10.1017/S1351324900000218

1331  Markman, A. B., & Makin, V. S. (1998). Referential communication and category

1332    acquisition. *Journal of Experimental Psychology. General*, *127*(4), 331–354.

1333    https://doi.org/10.1037/0096-3445.127.4.331

1334  Markman, E. (1989). *Categorization and Naming in Children*. Cambridge, MA: MIT

1335    Press.

1336  McKay, R., & Dennett, D. (2009). The evolution of misbelief. *Behavioral and Brain*

1337    *Sciences*, *32*(06), 493. https://doi.org/10.1017/S0140525X09990975

1338  McNicholas, P. (2016). Model-Based Clustering. *Journal of Classification*, *33*(3), 331–

1339    373. https://doi.org/10.1007/s00357-016-9211-9

1340  Mirza, M., Adams, R., Mathys, C., & Friston, K. (2016). Scene Construction, Visual

1341    Foraging, and Active Inference. *Frontiers in Computational Neuroscience*, *10*, 56.

1342    https://doi.org/10.3389/fncom.2016.00056

1343    Mirza, M. B., Adams, R. A., Mathys, C., & Friston, K. J. (2018). Human visual

1344        exploration reduces uncertainty about the sensed world. *PLOS ONE*, *13*(1),

1345        e0190429. https://doi.org/10.1371/journal.pone.0190429

1346    Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., …

1347        Hassabis, D. (2015). Human-level control through deep reinforcement learning.

1348        *Nature*, *518*(7540), 529–533. https://doi.org/10.1038/nature14236

1349    Osherson, D. N., & Smith, E. E. (1981). On the adequacy of prototype theory as a

1350        theory of concepts. *Cognition*, *9*(1), 35–58. https://doi.org/10.1016/0010-

1351        0277(81)90013-5

1352    Parr, T., & Friston, K. (2017). Working memory, attention, and salience in active

1353        inference. *Scientific Reports*, *7*(1), 14678. https://doi.org/10.1038/s41598-

1354        017-15249-0

1355    Parr, T., & Friston, K. (2018). The Anatomy of Inference: Generative Models and

1356        Brain Structure. *Frontiers in Computational Neuroscience*, *12*, 90.

1357        https://doi.org/10.3389/fncom.2018.00090

1358    Parr, T., Markovic, D., Kiebel, S., & Friston, K. (2019). Neuronal message passing

1359        using Mean-field, Bethe, and Marginal approximations. *Scientific Reports*, *9*(1),

1360        1889. https://doi.org/10.1038/s41598-018-38246-3

1361    Perfors, A., Tenenbaum, J., Griffiths, T., & Xu, F. (2011). A tutorial introduction to

1362        Bayesian models of cognitive development. *Cognition*, *120*(3), 302–321.

1363        https://doi.org/10.1016/j.cognition.2010.11.015

1364    Perogamvros, L., & Schwartz, S. (2012). The roles of the reward system in sleep and

1365        dreaming. *Neuroscience & Biobehavioral Reviews*, *36*(8), 1934–1951.

1366      https://doi.org/10.1016/J.NEUBIOREV.2012.05.010

1367  Pfeiffer, B., & Foster, D. (2013). Hippocampal place-cell sequences depict future

1368      paths to remembered goals. *Nature*, *497*(7447), 74–79.

1369      https://doi.org/10.1038/nature12112

1370  Ritter, S., Wang, J., Kurth-Nelson, Z., & Botvinick, M. (2018). Episodic Control as

1371      Meta-Reinforcement Learning. *BioRxiv*, 360537.

1372      https://doi.org/10.1101/360537

1373  Salakhutdinov, R., Tenenbaum, J. B., & Torralba, A. (2013). Learning with

1374      hierarchical-deep models. *IEEE Transactions on Pattern Analysis and Machine*

1375      *Intelligence*, *35*(8), 1958–1971. https://doi.org/10.1109/TPAMI.2012.269

1376  Sanborn, A. N., Griffiths, T. L., & Navarro, D. J. (2010). Rational Approximations to

1377      Rational Models: Alternative Algorithms for Category Learning. *Psychological*

1378      *Review*, *117*(4), 1144–1167. https://doi.org/10.1037/a0020511

1379  Schmidhuber, J. (2006). Developmental robotics, optimal artificial curiosity,

1380      creativity, music, and the fine arts. *Connection Science*, *18*(2), 173–187.

1381      https://doi.org/10.1080/09540090600768658

1382  Schwartenbeck, P., FitzGerald, T., Mathys, C., Dolan, R., & Friston, K. (2015). The

1383      Dopaminergic Midbrain Encodes the Expected Certainty about Desired

1384      Outcomes. *Cerebral Cortex*, *25*(10), 3434–3445.

1385      https://doi.org/10.1093/cercor/bhu159

1386  Schwartenbeck, P., & Friston, K. (2016). Computational Phenotyping in Psychiatry: A

1387      Worked Example. *ENeuro*, *3*(4), ENEURO.0049-0016.2016.

1388      https://doi.org/10.1523/ENEURO.0049-16.2016

1389  Sharot, T. (2011). The optimism bias. *Current Biology*, *21*(23), R941–R945.

1390  https://doi.org/10.1016/J.CUB.2011.10.030

1391  Stachenfeld, K. L., Botvinick, M. M., & Gershman, S. J. (2016). The hippocampus as a

1392  predictive map.

1393  Stickgold, R., Hobson, J., Fosse, R., & Fosse, M. (2001). Sleep, learning, and dreams:

1394  off-line memory reprocessing. *Science*, *294*(5544), 1052–1057.

1395  https://doi.org/10.1126/science.1063530

1396  Tervo, D. G. R., Tenenbaum, J. B., & Gershman, S. J. (2016). Toward the neural

1397  implementation of structure learning. *Current Opinion in Neurobiology*, *37*, 99–

1398  105. https://doi.org/10.1016/J.CONB.2016.01.014

1399  Tononi, G., & Cirelli, C. (2014). Sleep and the Price of Plasticity: From Synaptic and

1400  Cellular Homeostasis to Memory Consolidation and Integration. *Neuron*, *81*(1),

1401  12–34. https://doi.org/10.1016/J.NEURON.2013.12.025

1402  Walker, M., & Stickgold, R. (2010). Overnight alchemy: sleep-dependent memory

1403  evolution. *Nature Reviews. Neuroscience*, *11*(3), 218; author reply 218.

1404  https://doi.org/10.1038/nrn2762-c1

1405  Wang, J X, Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., …

1406  Botvinick, M. (2016). Learning to reinforcement learn. *ArXiv:1611.05763*.

1407  Wang, Jane X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z., …

1408  Botvinick, M. (2018). Prefrontal cortex as a meta-reinforcement learning

1409  system. *Nature Neuroscience*, *21*(6), 860–868.

1410  https://doi.org/10.1038/s41593-018-0147-8

1411  Ward, T. B. (1994). Structured Imagination: the Role of Category Structure in

1412      Exemplar Generation. *Cognitive Psychology*, *27*(1), 1–40.

1413      https://doi.org/10.1006/cogp.1994.1010

1414   Whittington, J. C. R., Muller, T. H., Mark, S., Barry, C., & Behrens, T. E. J. (2018).

1415      Generalisation of structural knowledge in the hippocampal-entorhinal system.

1416   Williams, J. J., & Lombrozo, T. (2010). The role of explanation in discovery and

1417      generalization: evidence from category learning. *Cognitive Science*, *34*(5), 776–

1418      806. https://doi.org/10.1111/j.1551-6709.2010.01113.x

1419   Wipf, D. P., & Rao, B. D. (2007). An Empirical Bayesian Strategy for Solving the

1420      Simultaneous Sparse Approximation Problem. *IEEE Transactions on Signal*

1421      *Processing*, *55*(7), 3704–3716. https://doi.org/10.1109/TSP.2007.894265

1422   Xu, F., & Tenenbaum, J. (2007a). Sensitivity to sampling in Bayesian word learning.

1423      *Developmental Science*, *10*(3), 288–297. https://doi.org/10.1111/j.1467-

1424      7687.2007.00590.x

1425   Xu, F., & Tenenbaum, J. (2007b). Word Learning as Bayesian Inference. *Psychological*

1426      *Review*, *114*(2), 245–272.

1427

1428