

1 **A Recurrent Neural Network Model for Flexible and Adaptive**
2 **Decision Making based on Sequence Learning**

3

4 Zhewei Zhang^{1,2}, Huzi Cheng¹, and Tianming Yang^{1*}

5

6 ¹ Institute of Neuroscience, Key Laboratory of Primate Neurobiology, CAS Center for
7 Excellence in Brain Science and Intelligence Technology, Shanghai Institutes for Biological
8 Sciences, Chinese Academy of Sciences

9 ² University of Chinese Academy of Sciences, Beijing 100049, China

10

11 *Address correspondence to:

12 Tianming Yang (tyang@ion.ac.cn)

13 Institute of Neuroscience

14 320 Yue Yang Rd.

15 System Neuroscience Building, Room 302

16 Shanghai, China 200031

17 Phone: +86-21-54921737, Fax: +86-21-54921735

18

19 **Abstract**

20

21 The brain makes flexible and adaptive responses in the complicated and
22 ever-changing environment for the organism's survival. To achieve this,
23 the brain needs to choose appropriate actions flexibly in response to
24 sensory inputs. Moreover, the brain also has to understand how its actions
25 affect future sensory inputs and what reward outcomes should be
26 expected, and adapts its behavior based on the actual outcomes. A
27 modeling approach that takes into account of the combined contingencies
28 between sensory inputs, actions, and reward outcomes may be the key to
29 understanding the underlying neural computation. Here, we train a
30 recurrent neural network model based on sequence learning to predict
31 future events based on the past event sequences that combine sensory,
32 action, and reward events. We use four exemplary tasks that have been
33 used in previous animal and human experiments to study different
34 aspects of decision making and learning. We first show that the model
35 reproduces the animals' choice and reaction time pattern in a probabilistic
36 reasoning task, and its units' activities mimics the classical findings of the
37 ramping pattern of the parietal neurons that reflects the evidence
38 accumulation process during decision making. We further demonstrate
39 that the model carries out Bayesian inference and may support meta-
40 cognition such as confidence with additional tasks. Finally, we show how
41 the network model achieves adaptive behavior with an approach distinct
42 from reinforcement learning. Our work pieces together many

43 experimental findings in decision making and reinforcement learning and
44 provides a unified framework for the flexible and adaptive behavior of the
45 brain.

46

47 **Introduction**

48

49 Consider a scenario in which a cheetah sneaks up on a deer until it starts the final dash to
50 catch it. Every move has to be calculated carefully based on the sensory inputs, for example,
51 the distance to the deer and the surrounding environment, and the past experience, for
52 example, the deer's speed. Furthermore, the cheetah should be able to predict how the deer
53 would respond to its move to allow more timely adjustment of its actions. This is an example
54 of a central function of the brain: to generate responses to maximize the gain based on a
55 continuous stream of sensory inputs and movements that reflect a complicated and volatile
56 world. The contingencies that the brains need to learn should take into account the sensory
57 inputs, the organism's own actions, and the resulting reward outcomes, all strung together into
58 a sequence.

59

60 A variety of theoretic models have been used to investigate certain components of this general
61 contingency problem. For example, in the perceptual decision-making and sensory-motor
62 transformation field, the research question often centers around how sensory information as
63 evidence is integrated and leads to specific actions. Theoretical models such as the drift-
64 diffusion model (Ratcliff, 1978; Stone, 1960) and attractor neural network models (X.-J.
65 Wang, 2001; Wong & Wang, 2006) have gained a lot of success in both modeling the behavior

66 and explaining the neuronal response patterns in the brain. Yet these studies often do not take
67 into account the volatile nature of the environment. In contrast, the conditioning and learning
68 field focus on the learning of the associations between reward or punishment outcomes and
69 other events in a volatile environment. The reinforcement learning (RL) framework has
70 become a standard approach to understand the behavior and the brain circuitry in this field
71 (Dayan & Niv, 2008; Schultz, Dayan, & Montague, 1997; Sutton & Barto, 2012). In RL, state
72 and action values are updated when reward feedbacks differ from expectations. The
73 framework provides an explanation on how animals gradually learn the values associated with
74 stimuli and actions and select actions accordingly. Yet, RL does not usually deal with the
75 problem of inferring states based on noisy sensory inputs. It is inefficient for learning
76 complicated decision-making tasks, especially those with a large state space. More
77 problematic is that in most studies RL assumes that states are observable and task structures
78 are known, which is not true in the real world. Biologically realistic neural network models of
79 RL that acquires task structures through learning have been proposed but they are limited to
80 rather simple structures (J. X. Wang et al., 2018; Zhang, Cheng, Lin, Nie, & Yang, 2018).

81

82 So far, there is a lack of a model general or flexible enough to cover different aspects of
83 decision making and learning. A more general framework is necessary for furthering our
84 understanding of the brain. Here, we propose a neural network model based on sequence

85 learning to model the flexible and adaptive behavior and to understand the underlying neural
86 mechanisms. The network takes inputs in the form of event sequences and is trained to make
87 predictions of future events with supervised learning. Importantly, the input event sequences
88 include not only the sensory events, but also the action events and the reward outcome events,
89 which allow the network to establish the contingencies between all three types of events. The
90 outputs of the network serve as a prediction of future events and the network is trained to
91 minimize the difference between the predictions and the future input sequences. The
92 predictions of action events are used to generate the model's actual responses, which are
93 assessed for the network's behavior performance. The predictions of sensory and reward
94 events are not directly exhibited in behavior, but they are part of the learning and have
95 important implications in modeling the brain.

96

97 To demonstrate how the network model works, we choose four exemplary behavior tasks that
98 have been previously used in animal and human studies, each covering a different aspect of
99 decision making and learning. The network is not constructed for any particular task. Instead,
100 it learns the task structure and task logic through training. The first task, a reaction time
101 version of probabilistic reasoning task, is one of the most complicated tasks in the literature
102 that provide us both behavior and single unit recording data to evaluate the model (Kira, Yang,
103 & Shadlen, 2015; Yang & Shadlen, 2007). We use it to demonstrate the capability of the

104 network to learn complex contingencies between events in long sequences and the similarity of
105 the units in the network to the previous neurophysiological findings. We further use a multi-
106 sensory integration task to show how the network generalizes Bayesian inference across
107 different sensory modalities (Gu, Angelaki, & DeAngelis, 2008). The third task, a post-
108 decision wagering task, is chosen to show that the model can support meta-cognition such as
109 confidence (Kiani & Shadlen, 2009). Last but not least, the final task is adapted from the two-
110 step task, first described to illustrate model-based RL (Daw, Gershman, Seymour, Dayan, &
111 Dolan, 2011). The task requires the learning to be extended across trials, which allows us to
112 demonstrate how the network accounts for adaptive behavior with a distinct approach from
113 RL.

114

115 Together, our results show that a network model that is simply trained to make predictions
116 based on event sequences may account for a large body of experimental findings from both
117 the decision making and the reinforcement learning field. Not only does the model reproduce
118 the animals' and humans' behavior in previous studies, the units in the network also show
119 response patterns resembling those of the neurons in the brain. The model makes testable
120 hypotheses on the neuronal response patterns and circuit structures in the brain and suggests
121 novel interpretations for some of the previous experimental work. These results suggest the

122 potential of our model to be a unified framework for decision making and learning that may
123 reveal the computational principle in the brain.

124 **Results**

125

126 **Network**

127

128 Our network model contains three layers: the input layer, the hidden layer based on gated
129 recurrent units (GRU), and the output layer (Fig 1a). The input layer contains units that carry
130 the information about the sensory, motor, and reward events on a timeline. Each unit's activity
131 can be a binary variable, representing the presence or the absence of the corresponding event,
132 or a continuous variable representing stimulus strength. The input units are fully connected
133 with the next layer. We use vector \vec{x}_j of length \mathfrak{N}_j to describe the activities of the input
134 layer units.

135

136 The hidden layer is the core of our model. It is adapted from the GRU network, which has
137 been shown to be suitable for learning sequences (Cho et al., 2014; Chung, Gulcehre, Cho, &
138 Bengio, 2014). There are $\mathfrak{N}_L=128$ units in the hidden layer. Each gated recurrent unit's
139 activity is a nonlinear combination of both the inputs and its activity at the previous time step,
140 which are regulated by an update gate and a reset gate. The nonlinear reset gate has been
141 shown to be critical for the network to learn long contingencies across long temporal
142 sequences (Greff, Srivastava, & Koutník, 2016). The state vector L_j of length \mathfrak{N}_L describes

143 the responses of these units at time t , which is updated with the information from the input

144 layer as follows:

145

146 $J_L = \text{Sigmoid}(\gamma_b + \gamma_L L_{t-1} + b_L),$ (1)

147 $R_L = \text{Sigmoid}(\lambda_b + \lambda_L L_{t-1} + b_R),$ (2)

148 $C_t = (1 - J_L) \cdot L_{t-1} + J_L \cdot \tanh(\gamma_b + \gamma_L (R_L \cdot L_{t-1}) + b_C),$ (3)

149 $L_t = [C_t],$ (4)

150

151 where J_L and R_L represent the update gate vector and the reset gate vector, respectively;

152 γ_b γ_L and γ_C are the input connection weight matrices for the update gates, the reset gates,

153 and the gated units; b_L b_R and b_C are the recurrent connection weight matrices for the update

154 gates, the reset gates, and the gated units; b_{J_L} b_{R_L} and b_{C_t} are the bias vectors for the update

155 gates, the reset gates, and the gated units; and \bullet indicates the element-wise multiplication. The

156 rectified-linear function $[\cdot]$ keeps the responses L_t non-negative. The initial value of the

157 state vector L_t at the beginning of each trial is reset to zero.

158

159 The hidden layer units project to the output layer with full connections. The output layer is

160 composed of an array of units that mirror the input layer units, representing the network's

161 predictions of the corresponding sensory events, action events and reward events for the next

162 time step. The activity of the output layer units \mathbf{s}_t (a vector of length $|\mathcal{A}|$) is:

163

164 $\mathbf{s}_t = \text{Sigmoid}(\mathbf{x}_t \mathbf{L}_o),$ (5)

165

166 where \mathbf{x}_t is the output connection weight matrix. The prediction y is a function of response

167 \mathbf{s}_t :

168

169 $\mathbf{r}_t = \frac{\mathbf{s}_t \mathbf{A}}{\mathbf{s}_t \mathbf{B}}$ (6)

170

171 The network predicts the corresponding event would happen in the next time step when \mathbf{r}_t^z

172 is 1. If \mathbf{r}_t^z corresponds to an action, the probability of the action being carried out at the

173 next time step is a softmax function of \mathbf{s}_t :

174 $\mathbf{p}_t @ \mathbf{a}_t | \mathbf{s}_t^z = \mathbf{e}^{\mathbf{s}_t^z / \tau} / \sum \mathbf{e}^{\mathbf{s}_t^z / \tau},$ (7)

175 where τ is the temperature. The chosen action is evaluated for the network's task

176 performance for testing purposes.

177

178 Fig 1a illustrates the network structure with the inputs and outputs corresponding to the
179 Task 1 that we discuss below. For the other tasks, it is only necessary to change the input
180 and output units to reflect the relevant events.

181

182 The goal of the training is for the network to learn to predict events and generate sequences
183 that lead to a reward. The loss function L is defined as the sum of mean squared error between
184 elements in the output \hat{s}_t and actual event sequence s_{t+4} for all time points t .

185

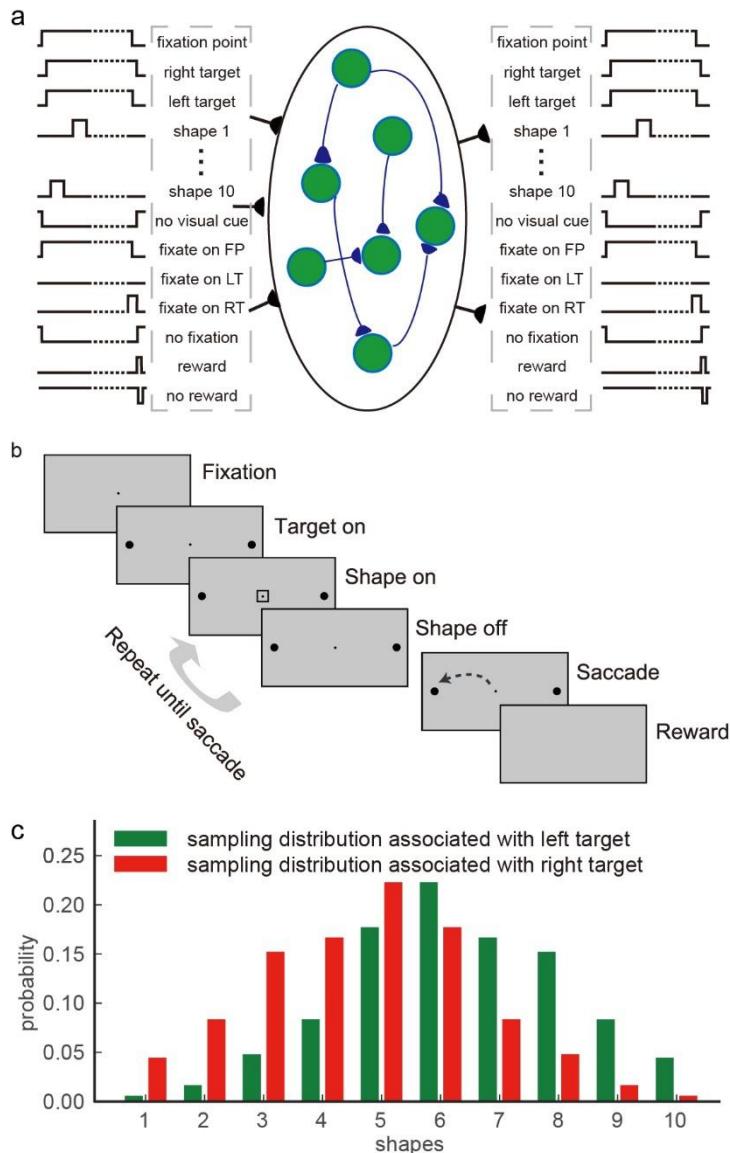
$$L = \frac{1}{T} \sum_{t=1}^{T-4} (\hat{s}_t - s_{t+4})^2 \quad (8)$$

187

188 where T is the length of each trial. The parameter vector θ (including θ_1 and θ_2) of the
189 network is updated with the Adam (Kingma & Ba, 2014) and the gradient clipping is applied
190 to avoid exploding gradients as follows:

191

$$\theta_t @ \frac{\partial L}{\partial \theta} \text{clip}(\theta_t, \theta_{t-1}, \theta_{t-1} + \frac{3.58}{\sqrt{\theta_t}}), \quad (9)$$



193

194 Figure 1. **a.** The model diagram. The network has three layers: the input layer, the gated
 195 recurrent unit layer, and the output layer. The input layer receives the input sequences of
 196 sensory events, action events, and reward events. The GRU layer has 128 gated recurrent
 197 units. The output layer units mirror the input layer units and represent the prediction of the
 198 future events. The diagram illustrates the particular input and output units for the Task 1.
 199 **b.** Task 1: the reaction-time version of probabilistic reasoning task. The subject fixates at
 200 a central point and views a sequence of shapes to make a decision by moving the eyes
 201 toward one of the two choice targets in the peripheral. Each shape confers information
 202 regarding to which target will be rewarded. An optimal strategy is to integrate the

203 information and make a choice when the integrated information hits a bound. **c.** The
204 sampling distributions. Shapes are sampled from the sampling distribution associated with
205 the correct target in each trial.

206

207 **Task 1: Probabilistic reasoning task**

208

209 We train our model with a reaction-time version of the probabilistic reasoning task that was
210 used to study the neural mechanism of sequential decision making (Kira et al., 2015). The task
211 is illustrated in Fig 1b. In this task, a subject has to make decisions between a pair of eye
212 movement targets. Initially, the subject needs to fixate on a central point on a computer screen.
213 Then a stream of shapes appears sequentially near the fixation point. There are totally 10
214 possible shapes. Each shape conveys information on the correct answer. The subject needs to
215 integrate information of the shapes to form a decision and move the eyes to look at the choice
216 target whenever ready.

217

218 The contingency between the shapes and the targets is described by two distributions. Each
219 target is associated with a distribution of the appearance probability of the shapes in the
220 sequence (Fig 1c). In each trial, the computer randomly picks the correct target. Each shape in
221 the sequence is independently sampled from the distribution associated with the correct target.
222 Because the likelihoods of observing a particular shape under the two distributions are

223 different, each shape provides information on which target is correct. It has been shown that
224 the sequential probability ratio test (SPRT) is an optimal strategy to solve the task in the sense
225 that it requires the least number of observations to achieve a desired performance (Wald &
226 Wolfowitz, 1948). In the SPRT, one needs to accumulate the log likelihood ratio (logLR)
227 associated with each piece of evidence, which is the log ratio between the conditional
228 probabilities of observing the evidence given the two testing hypotheses. In our task, the
229 logLRs associated with each shape range from -0.9 to 0.9 (base 10). We define positive and
230 negative logLRs as evidence supporting the left and the right target, respectively.

231

232 The task has several attractive features for our modeling purposes. First of all, the task
233 features a sophisticated statistical structure containing the shapes, the choice, and the reward.
234 As the shapes appear one by one, an ideal observer not only gains information on what would
235 be the appropriate choice, but also can deduce how likely a particular shape will appear next
236 and how likely a reward can be expected. Second, this is a reaction time task in which choices
237 have to be made at appropriate time to achieve a certain tradeoff between speed and accuracy.
238 This allows us to demonstrate the flexibility of our model for learning sequences of variable
239 length. Last but not least, the task is one of the most complicated tasks that have been used in
240 animal studies with both behavior and neural data available. We not only can compare the

241 behavior between the model and the animals, but also can look into the network and study the
242 network units' activities and compare them against the experimental findings.

243

244 The training dataset

245

246 We train the network with task event sequences created with simulated trials in which choices
247 are generated with a drift-diffusion model with collapsing bounds. In each trial, a correct
248 target is randomly determined, and the shapes are generated with its associated distribution
249 (Fig 1c). A choice is triggered when the accumulated logLR reaches to either of the two
250 opposite collapsing bounds. The bounds start at ± 1.5 and linearly collapse toward 0 at the rate
251 of 0.1 per shape epoch. The left target is chosen if the positive bound is reached, and the right
252 target chosen if the negative bound is reached. If the choice matches the pre-determined
253 correct target, a reward is given, and the corresponding sequence is included in the training
254 dataset.

255

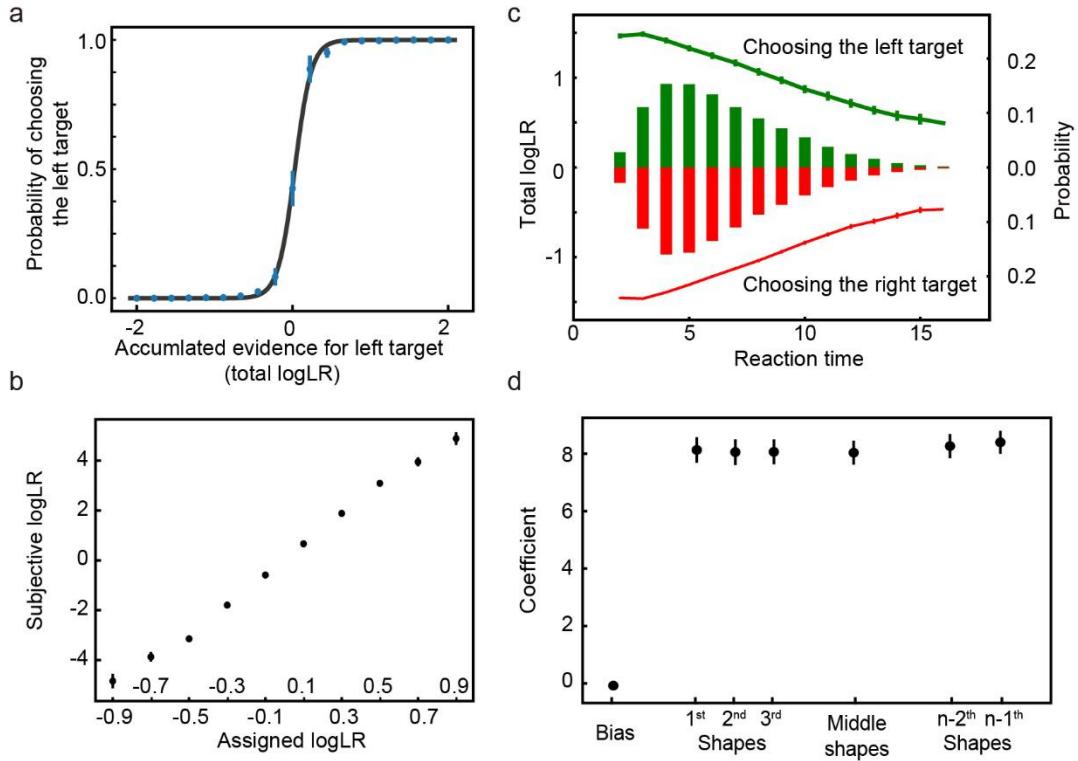
256 The results presented below are based on 20 simulation runs using a training set generated
257 with 7.5×43^8 simulated trial sequences.

258

259 Behavior Performance

260

261 After training, the network performs the task well. When the total logLR associated with the
262 shape sequence is positive, the network tends to choose the left target, and when the total
263 logLR is negative, it tends to choose the right target (Fig 2a). A logistic regression reveals that
264 the logLR assigned to each shape correlates well with its leverage on the choice (Fig 2b). The
265 mean reaction time, quantified as the number of shapes that the model uses for decision, is
266 6.36 ± 0.04 (mean \pm s.e.m.). Both the distribution of the reaction time and the total logLR at the
267 time of choice suggest that the model behavior is consistent with the DDM with collapsing
268 bounds (Fig 2c). We use another logistic regression to examine the effects of shape order on
269 the choice. The first shapes in the sequence exert similar leverages on the choice to the rest of
270 the shapes, except for the last shape, suggesting all shapes except the last are used equally in
271 decision making (Fig 2d). The last shape's sign is almost always consistent with the choice ($>$
272 99% of all trials). This is also consistent with the DDM, in which the last shape brings the total
273 logLR over the bound. Overall, the network performance resembles, although understandably
274 better than, the behavior of the macaque monkeys trained with a very similar task (Kira et al.,
275 2015).



276

Figure 2. **a.** The psychometric curve. The model more often chooses the target supported by the accumulated evidence. The black curve is the fitting curve from the logistic regression. **b.** The leverage of each shape on choice revealed by the logistic regression is consistent with its assigned logLR. **c.** Reaction time. The bars show the distribution of reaction time, quantified by the number of observed shapes (right y-axis). Green and red indicate the left and right choices, respectively. The lines indicate the mean total logLR (left y-axis) at the time of decision, grouped by reaction time. Trials with only 1 shape or more than 16 shapes comprise less than 0.1% of the total trials and are excluded from the plot. **d.** The leverage of the first 3, the second and third from the last, and the middle shapes on the choice. Only trials with more than 6 shapes are included in the analysis. No significant differences are found between any pair of the coefficients of shape regressors (two-tailed t test with Bonferroni correction). The error bars in all panels indicate S.E. across runs. Some error bars are smaller than the data points and not visible.

290

291 The good performance of the network is not because the trials sequences in the testing dataset
292 overlap with those in the training dataset significantly, which is not true because the possible
293 number of trial sequences in this task is astronomical and each shape sequence in the training
294 and testing dataset is randomly generated. The network is able to generalize beyond the
295 training dataset. With a training dataset that contains only 1000 unique trial sequences, the
296 network can still achieve a good performance (Supplementary Figure 1). Therefore, the
297 learning of the statistical structure of the task is likely the reason for the model's performance.

298

299 Network analyses – Evidence and Choice encoding

300

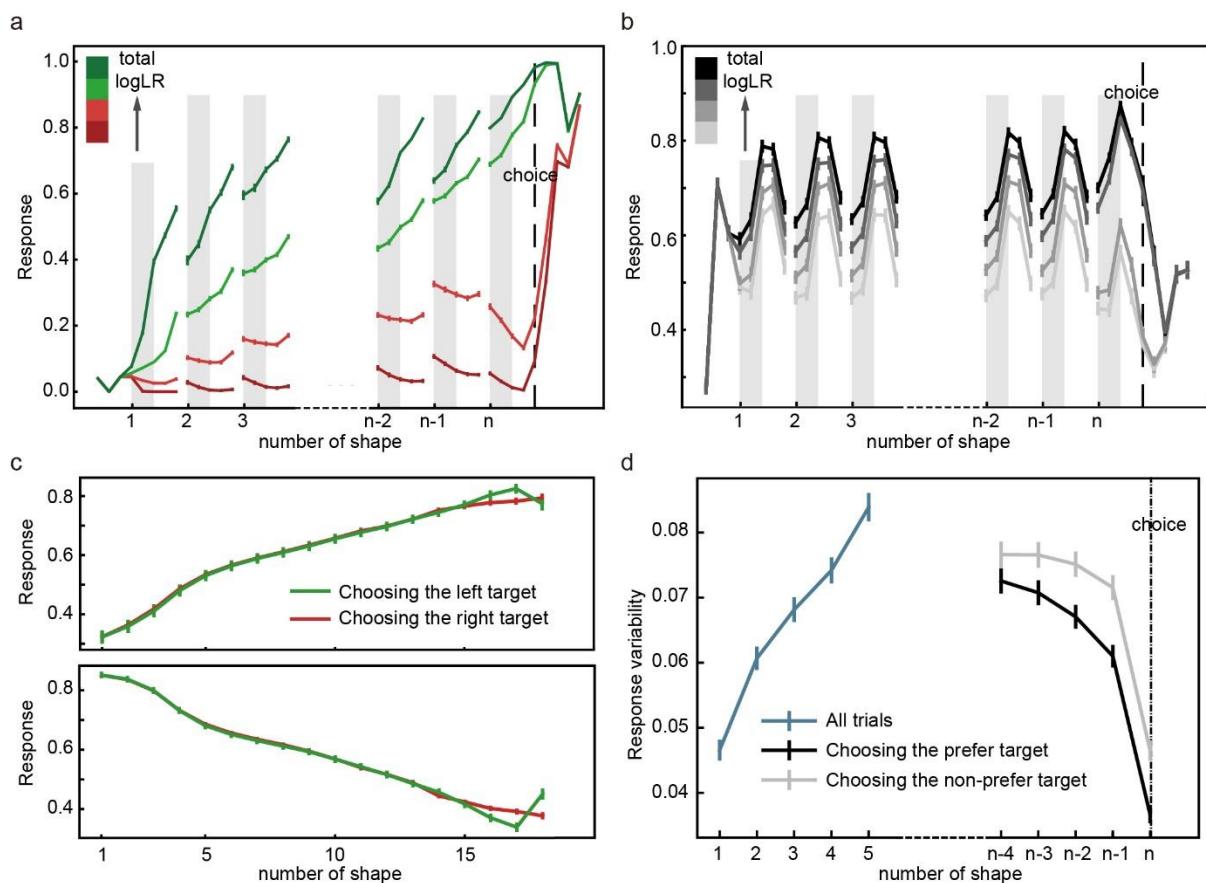
301 Next, we examine how the evidence and the choice are encoded in the network. The example
302 unit in Fig 3a shows a classical ramping-up activity pattern that has been reported in neurons
303 from the prefrontal cortex (Kim & Shadlen, 1999), the parietal cortex (Kira et al., 2015;
304 Roitman & Shadlen, 2002; Yang & Shadlen, 2007), and the striatum (Ding & Gold, 2010). Its
305 activity increases when the total logLR grows to support its preferred choice and decreases
306 when the total logLR is against its preferred choice. The responses converge around the time
307 when its preferred choice is chosen. The population analysis using all the units that are
308 selective to the total logLR finds the same pattern (Fig 3b).

309

310 In addition to the units that accumulate evidence and reflect the decision making process, we
 311 also find units that show a ramping-up or ramping-down activity pattern independent of the
 312 choice (Fig 3c). Their activities indicate the passage of time and may be interpreted as an
 313 urgency signal. Neurons with a similar activity pattern have been reported in the global
 314 pallidus (Thura & Cisek, 2017).

315

316



317

318 Figure 3. **a.** An example unit that prefers the left target. Its activity increases when the
 319 evidence supporting the left target grows and decreases when it drops. The unit's responses
 320 converge when the network chooses its preferred target. The trials are grouped into

321 quartiles by the accumulated evidence in each epoch, which are indicated with color. The
322 error bars indicate the S.E. across trials. **b.** Population responses of the units that are
323 selective to the total logLR. The trials are grouped based on the accumulated evidence
324 supporting each unit's preferred target in each epoch. The error bars in panels **b**, **c** and **d**
325 indicate the S.E. across units. **c.** Urgency units. Their activities ramp up (upper panel) or
326 down (lower panel) regardless of the choice. **d.** Network unit response variability. The
327 neurons' response variability increases initially (blue curve) but decreases before the
328 choice, more so when the preferred target is chosen (black) than when the non-preferred
329 target is chosen (grey). Only the trials with more than five shapes are included in panel **a**,
330 **b** and **d**.

331

332 We use a linear regression to quantify the selectivity of each units (see Methods) to three
333 important parameters for this task: the accumulated evidence, quantified with the total logLR,
334 the choice outcome, and the urgency. During the shape presentation period, a large portion of
335 units is found to encode the accumulated evidence ($N = 63.00 \pm 2.49$), the choice outcome (N
336 $= 50.60 \pm 2.03$) and the urgency ($N= 86.50 \pm 1.10$).

337

338 Finally, we calculate the response variability of the units in the hidden layer to study the
339 dynamics of the population responses of the network (Fig 3d). This is analogous to the
340 variance CE that Churchland and colleagues studied previously in the LIP neurons
341 (Churchland et al., 2011). Our analyses are more straightforward as the units in our model do
342 not have intrinsic Poisson noise. The response variability here is simply calculated as the

343 standard deviation of the units' responses across trials. We find that the response variability
344 increases initially (linear regression, $k = 0.0088$, $p < 0.001$) as more shapes are presented and
345 the evidence is accumulated. However, when we align the trials to the choice, a different
346 pattern emerges. In the trials in which the preferred target is chosen, the response variability
347 decreases before the choice (linear regression, $k = -0.0038$, $p < 0.001$; Spearman's rank
348 correlation, $r = -0.045$, $p < 0.01$) and reaches the minimum around the time of choice. In
349 contrast, in the trials in which the non-preferred target is chosen, the response variability is
350 significantly higher and its decrease is not significant until the last shape (linear regression, k
351 $= -0.0017$, $p = 0.07$; Spearman's rank correlation, $r = -0.017$, $p = 0.28$). The overall pattern is
352 very similar to that of the LIP neurons reported previously (Churchland et al., 2011).

353

354 Network analyses – When and Which

355

356 The balance between the speed and accuracy is an important aspect of decision making. The
357 control of the speed and accuracy balance has been suggested to be exerted through the same
358 neurons that accumulate evidence (Ding & Gold, 2012; Hanks, Kiani, & Shadlen, 2014; Heitz
359 & Schall, 2012) or a distinct populations of neurons that reflect only the speed but not the
360 choice (Thura & Cisek, 2017). We analyze the units' activities and connectivities in the model
361 to find out which is the case here in our model.

362

363 The network has three output units related to the choice. They include a unit for the fixation
364 (FP) and two units for the saccades toward the left (LT) and the right (RT) target, respectively.
365 We examine the connection weights between the hidden layer unit and the output units to
366 understand how hidden layer units drive the choices. For each hidden layer unit, we estimate
367 how it contributes to the outputs by defining two indices, I_{when} and I_{which} :

368

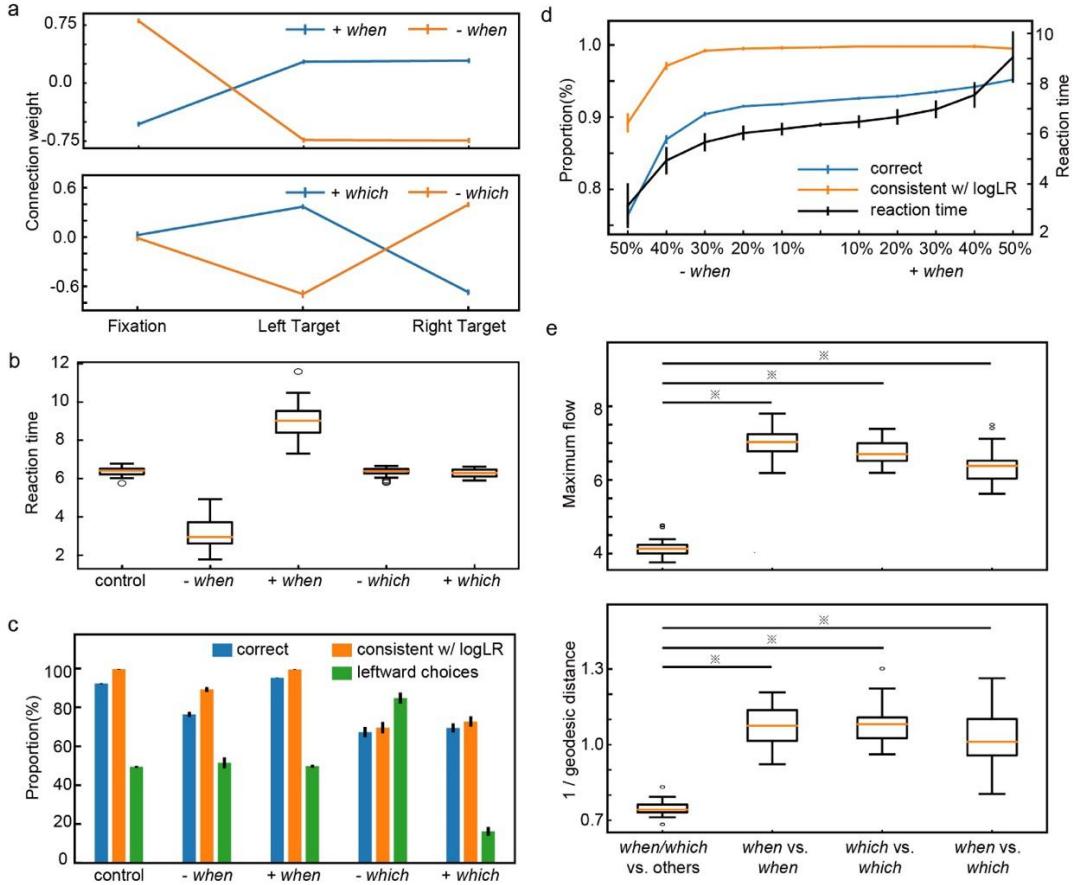
$$I_{when} = +\frac{L_{LT}}{L_{LT} + L_{RT}} - \frac{L_{RT}}{L_{LT} + L_{RT}} \quad (10)$$

$$I_{which} = \frac{L_{LT}}{L_{LT} + L_{RT}}, \quad (11)$$

371

372 where L_{LT} , L_{RT} and L_{FP} are the connection weights between the hidden layer unit
373 i and the output layer units LT, RT and FP, respectively. Units with large positive I_{when}
374 promote the saccades over fixation, regardless of the direction, thus affect the reaction time.
375 In contrast, units with large I_{which} bias the saccade direction toward their preferred direction.

376



377

378 Figure 4. **a**. The connection weights between the eye movement output units and the *when*
 379 units (upper panel) and the *which* units (lower panel). **b**. The effect of lesions to the *when*
 380 and *which* units on the reaction time. **c**. The effect of lesions to the *when* and *which* units
 381 on the choice. The blue bars indicate the proportion of correct trials. The orange bars
 382 indicate the proportion of trials in which the choice is consistent with the sign of the
 383 accumulated evidence at the time of choice. The green bars indicate the percentage of
 384 trials in which the model chooses the left target. **d**. Speed-accuracy tradeoff. We suppress
 385 the output of a different portion of *+when/-when* units (see Methods). As more *+when*
 386 units' outcome is suppressed, the model's reaction time (black curve, right y-axis)
 387 increases along with the accuracy (blue curve, left y-axis). However, the proportion of
 388 trials in which the choices are consistent with the evidence (orange curve, left y-axis) stays
 389 the same except for the extreme cases. **e**. The maximum flow (upper panel) and the
 390 inverse of the geodesic distance (lower panel) between different unit groups. The smaller

maximum flow and the larger geodesic distance between *when/which* units and other units suggest the relatively tight connections between the *when* and *which* units. * indicates significant difference ($p < 0.05$, Two tailed t-test with Bonferroni correction). The error bars in all panels indicate the S.E. across runs.

With these indices defined, we select two groups of hidden layer units: the *when* units and the *which* units. Within each group, we further divide them into a positive subgroup and a negative subgroup, depending on their signs of I_{when} or I_{which} . The +*when* units are selected with the following procedure. First, we sort all units with positive I_{when} values by their I_{when} values in the descending order. Then, we calculate the accumulative I_{when} along this axis and select the top units that together contribute more than 50% of the sum of the positive I_{when} . These units are defined as the +*when* units (10.70 ± 0.26 out of 50.60 ± 0.89 units with positive I_{when}). They have larger connection weights to the two saccade output units than to the fixation output unit (Fig 4a). We use similar procedures to select the -*when*, the +*which*, and the -*which* units ($n = 13.80 \pm 0.27$ out of 77.40 ± 0.89 units with negative I_{when} , 9.80 ± 0.30 out of 64.00 ± 1.25 units with positive I_{which} , and 9.50 ± 0.28 out of 64.00 ± 1.25 units with negative I_{which} , respectively). The -*when* units have smaller connection weights to the saccade units than to the fixation unit, the +*which* units have larger connection weights to the left saccade than to the right saccade output unit, and the -*which* units have smaller connection weights to the left saccade than to the right saccade output unit (Fig 4a). With these selection criteria, the *when* units are supposed to contribute more to the decision of when a response should be made,

411 while the *which* units contribute more to the decision of which target should be chosen.

412 Together, they constitute $34.22 \pm 0.52\%$ of all the units in the hidden layer.

413

414 Next, we examine how the \pm *when* and \pm *which* units affect the network's behavior. The

415 causality link can be demonstrated by simulated lesions that selectively deactivate the

416 connections of each group of units to the output layer. With only the outputs inactivated, the

417 hidden layer itself is not disturbed. When we deactivate the $+$ *when* units, the RTs become

418 longer while the accuracy remains intact (Fig 4bc). In contrast, when we deactivate the output

419 of the $-$ *when* units, the network's RTs become smaller (Fig 4b). Although the network

420 performance has an apparent drop, the network's choice still accurately reflects the

421 accumulated logLR, suggesting the performance drop is largely due to the fact that the

422 network has to work with a smaller total logLR due to the shorter RT (Fig 4c). In comparison,

423 when *which* units are manipulated, only the choice accuracy is affected but the RT remains the

424 same (Fig 4bc). More specifically, inactivating $+$ *which* units leads to a bias toward the right

425 target, while inactivating $-$ *which* units leads to a bias toward the left target (Fig 4c). These

426 results suggest there are two distinct populations of hidden layer units contributing to the

427 choice and the reaction time.

428

429 The way how *when* units affect the accuracy and reaction time suggests a possible mechanism
430 to modulate the speed-accuracy tradeoff. We demonstrate this by suppressing different
431 amount of *when* units' outputs (see Methods for details). As more +*when* units' output is
432 suppressed, the reaction time becomes larger (Spearman's rank correlation, p<0.001) and the
433 accuracy becomes higher (Spearman's rank correlation, p<0.001) (Fig 4d). Importantly, the
434 choices are still consistent with the accumulated evidence except for the most extreme cases.

435

436 The distinct functional roles of the *when* and *which* units may suggest they have different
437 connection patterns. We use the connection matrix of the hidden layer units to construct a
438 weighted directed graph and calculate the geodesic distance and the maximum flow between
439 the hidden layer units (see Methods). To account for the units that are not connected to each
440 other, we use the average maximum flows and inverse of geodesic distance in the analysis,
441 which is 0 for the unit pairs that are not connected. Interestingly, we find the inverses of the
442 geodesic distance between the *when* units and the *which* units are significantly larger than that
443 between the *when/which* units and others. The maximum flows between the *when* units and
444 the *which* units are also significantly higher than the network average (Fig 4e). These results
445 suggest that the *when* and the *which* units belong to a tightly connected sub-network within the
446 hidden layer and are not topologically separated.

447

448 Network analyses – Predictive Coding

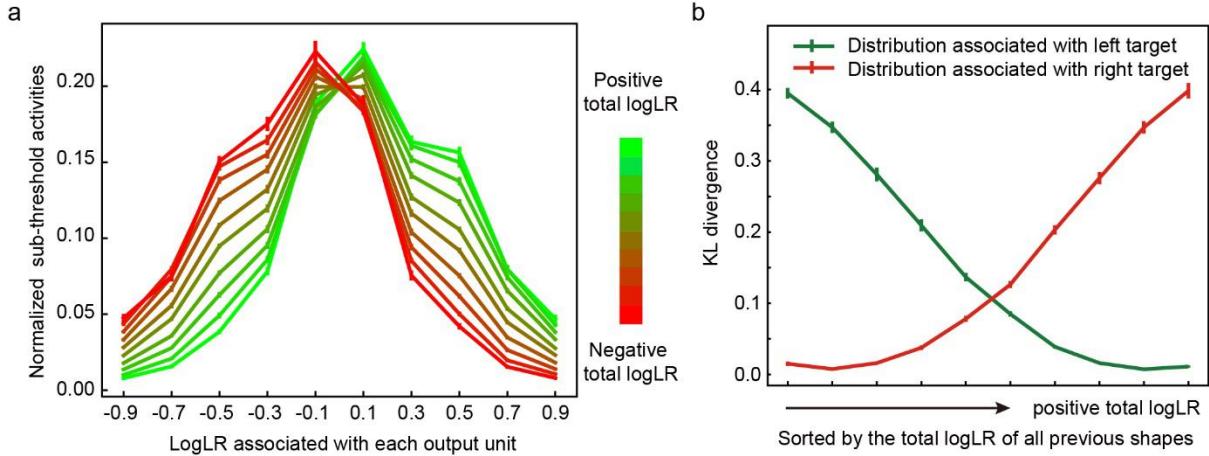
449

450 So far, we have shown that our network is able to generate appropriate choices based on shape
451 sequences. However, the network is constructed and trained in the way that the sensory events
452 are treated exactly the same as the action events. Therefore, the activities of the output units
453 that correspond to the sensory events should provide predictions on the sensory events. In
454 other words, our network may also serve as a generative model for the predictive coding
455 framework.

456

457 In the probabilistic reasoning task, each choice is associated with its respective shape
458 distribution. Therefore, as one accumulates information regarding to the choice, the
459 probability distribution of the upcoming shapes can also be inferred. We plot the mean
460 subthreshold activity, o_r , of the ten shape output units in Figure 5a at the time point right
461 before the onset of each shape. When the current evidence is in favor of the left target, the
462 activities of the output layer shape units resemble the probabilities from which the shapes are
463 drawn when the left target is the correct target. When the evidence is in favor of the right
464 target, the activity pattern of the output layer shape units shifts and becomes more similar to
465 the sampling distribution associated with the right target (Fig 5a).

466



467 Figure 5. **a.** The normalized subthreshold activities of 10 shape output units. We calculate
 468 the 10 shape units' activities at the time step immediately before each shape onset for all
 469 epochs in all trials in the test dataset (96555 trials and 583445 epochs from 20 runs),
 470 normalized by dividing each unit's activity by the sum of activities of all 10 shape output
 471 units. Data are divided into 10 group by the total logLR before the shape onset, which is
 472 indicated by the color. **b.** The Kullback-Leibler (KL) divergence between the normalized
 473 subthreshold activities (as shown in the Fig. 5a) and the sampling distributions (shown in
 474 the Fig. 1c). Data are grouped by the total logLR. The error bars indicate the S.E. across
 475 runs.
 476

477
 478 We quantify the similarity between the activity patterns of the output layer shape units and the
 479 sampling distributions by calculating the Kullback-Leibler (KL) divergence (Fig 5b). The KL
 480 divergence between the output layer unit activities and the sampling distribution associated
 481 with the left target decreases as the total logLR supporting the left target gets larger and grows
 482 when the total logLR is smaller. The opposite trend is observed on the KL divergence between
 483 the output layer unit activities and the sampling distribution associated with the right target.

484 These results suggest the activities of the output layer shape units encode the probability
485 distribution of the next shape based on the current accumulated evidence.

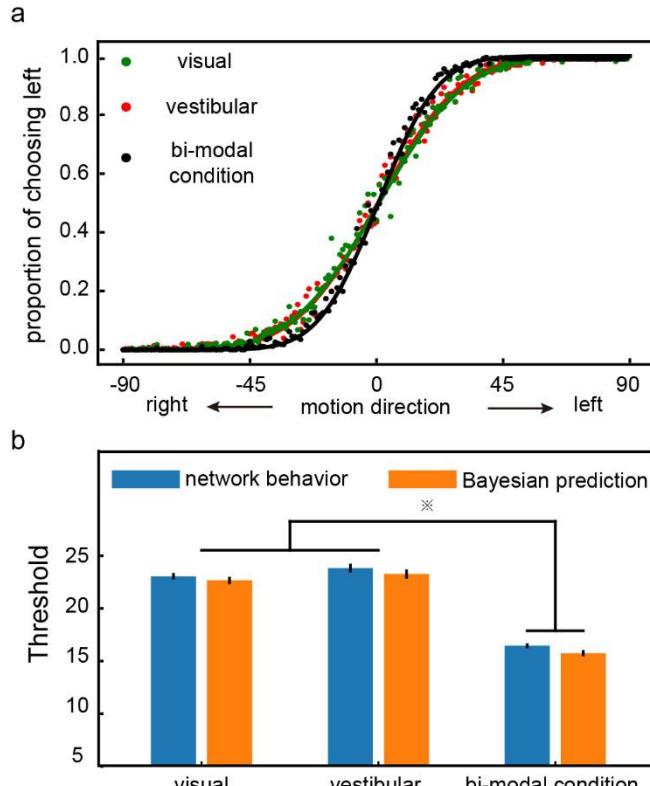
486

487 **Task 2: Multi-sensory integration task**

488

489 So far, we have shown that the network is able to perform the probabilistic reasoning task.
490 The decision making process in the task may also be interpreted as Bayesian inference in
491 which each shape updates the prior information. To more directly illustrate whether the
492 network performs Bayesian inference, we test the model with a multi-sensory integration task
493 adapted from Gu et al., 2008 (Gu et al., 2008) and compare it with the animals' behavior. In
494 this task, the subject has to discriminate whether the heading direction is toward the left or the
495 right basing on either the visual input, the vestibular input, or both. Gu and colleagues showed
496 that the animals were able to use information combined from different modalities optimally
497 and reached a threshold that matched predictions based on Bayesian inference.

498



499

500 Figure 6. **a.** The psychometric curve. The model is first trained with the single modal
 501 conditions and then tested with both the single modal (green: visual, orange: vestibular)
 502 and the bi-modal (black) conditions. Each data point represents the proportion of the left
 503 choice at a given motion direction condition. The model shows a steeper psychometric
 504 curve for the bi-modal condition, indicating a better performance. **b.** The performance
 505 thresholds. The blue bars are the model's thresholds under the single modal and the bi-
 506 modal conditions, compared against the thresholds calculated with the optimal Bayesian
 507 inference (orange). The threshold under the bi-modal condition is significantly smaller
 508 than the thresholds under either single modal condition. The differences between the
 509 thresholds of the network and the thresholds calculated with Bayesian inference are not
 510 significant. (two-tailed t-test with Bonferroni correction, p value threshold = 0.05)

511

512 To model this task, we assume that the input layer units have Gaussian heading-direction
 513 tuning curves with Poisson noise. They are divided into two groups of 8 units, corresponding

514 to the visual and vestibular inputs, respectively. Critically, the network is trained with trials in
515 which inputs are from only one single sensory modality and tested with trials in which inputs
516 are still limited to single modalities and trials in which inputs from both modalities are
517 available. We find that the network can generalize the training to bi-modal inputs. It combines
518 the information from visual and vestibular inputs and achieves a higher accuracy without
519 further training (Fig 6a). We fit the network model's choices in both the uni-modal and bi-
520 modal conditions with cumulative Gaussian functions and define network's performance
521 thresholds as the standard deviation of the best-fitting function. The threshold under the
522 combined condition is significantly smaller than the thresholds under either single modality
523 conditions (Fig 6b). Importantly, it matches the predictions of Bayesian inference. These
524 results suggest that the network model is able to perform Bayesian inference without being
525 explicitly trained to do so.

526

527 [Task 3: Confidence / Post-decision wagering task](#)

528

529 It has been argued that the same neurons that represent the decision variable during decision
530 making may also support meta-cognition such as confidence. We test how our model performs
531 a post-decision wagering task (Kiani & Shadlen, 2009). In this task, the model needs to make
532 decisions about the movement direction of a random dot motion stimulus. The task difficulty

533 is controlled by the proportion of dots moving coherently and the duration of the random dot
534 stimulus. A reward is delivered if the decision matches the direction of the coherently moving
535 dots. In half of the trials, a third target (sure target) appears after the motion viewing period.
536 Choosing the sure target leads to a smaller but certain reward.

537

538 We set up the appropriate inputs and output units for our network model to learn the task.
539 There are 10 input units for the motion stimulus, among which 5 prefer the leftward and the
540 other 5 prefer the rightward motion. Their activities represent the motion strength of the
541 random dots. Independent noises are added to the activities. Two additional input units
542 indicate the left and the right choice target. Finally, there is an input unit that indicates the
543 presence of the sure target. Again, the output units are a mirror copy of the input units, and
544 the activities of the output units corresponding to the fixations on the choice targets and the
545 sure target are used for generating the decisions.

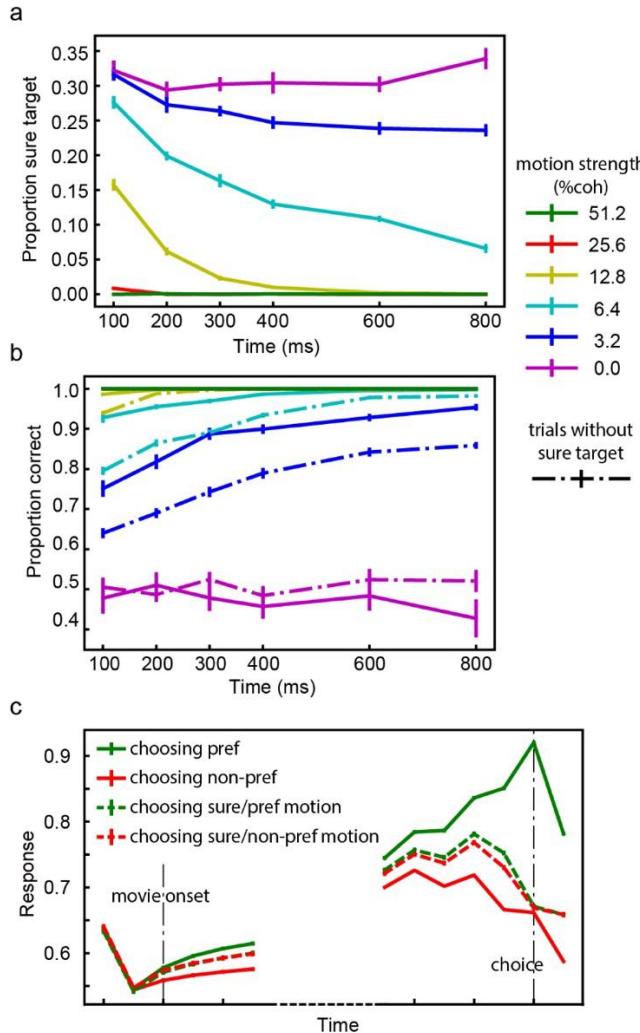
546

547 The network exhibits a similar choice pattern to the monkeys' behavior (Kiani & Shadlen,
548 2009). When the motion strength is weaker and the stimulus duration is shorter, the task
549 difficulty is higher and the network chooses the sure target more frequently (Fig 7a). In
550 addition, because the model may opt out for the sure target when its confidence about the
551 motion direction is low, the overall accuracy becomes higher in trials in which the sure target

552 is available but the model chooses to indicate the motion direction than in trials without the
553 sure target (Fig 7b). The model's behavior is consistent with the monkeys' strategy.

554

555 We further look at the response patterns of the units in the hidden layer. In particular, we want
556 to study the units that show choice selectivity during the delay period before the choice. This
557 selection criterium is similar to what was employed in the experimental study (Kiani and
558 Shadlen, 2009, see Methods). The responses of these units are modulated during the motion
559 viewing period (Fig 7c). More importantly, in the trials when the network chooses the sure
560 target, their responses reach an intermediate level between the activity levels in trials when the
561 targets of the preferred and non-preferred directions are chosen. These units behave like the
562 neurons recorded from the LIP area and may provide a basis for the post-decision wagering
563 (Kiani & Shadlen, 2009).



564

565 Figure 7. **a.** The proportion of trials in which the model chooses the sure target. The color
 566 indicates the motion strength. The frequency of the sure target choices decreases with the
 567 motion strength and the duration of motion viewing. **b.** The accuracy. Solid lines are the
 568 trials with the sure target and dashed lines are the trials without the sure target. The correct
 569 rate is higher when the sure target is provided. The error bars in panels **a** and **b** are S.E.
 570 across runs. **c.** Activities of choice-selective units. The responses are aligned to the movie
 571 onset (left) and the choice (right). The color of each line denotes the choice and motion
 572 direction. The dashed lines are trials in which the model chooses the sure targets. The units
 573 have intermediate activity level in trials that the sure target is chosen. The error bars are
 574 S.E. across units.

575

576

577 **Task 4: Two step task**

578

579 In the last experiment, we extend our model to a learning task. Learning requires the model to
580 incorporate trial history information for future decisions. By piecing events across multiple
581 trials together, our network model should be able to infer contingencies across trials and
582 achieve adaptive behavior. Notably, because the sensory, motor, and reward events are all
583 included in our model, it is straightforward to achieve across-trial learning with our model.

584

585 To demonstrate that the network adapts well in changing environment, we use an adapted
586 version of the origin two-step task (Akam, Costa, & Dayan, 2015; Daw et al., 2011). In this
587 task, the agent has to choose between two options A1 and A2, each leading to one of the two
588 intermediate outcomes, B1 and B2, with different but fixed transition probabilities. In
589 particular, A1 is more likely followed by B1 (common) than B2 (rare) and A2 is more likely
590 to be followed by B2 (common) than B1 (rare). Finally, B1 and B2 are each associated with a
591 probabilistic reward (0.8 vs 0.2). The reward contingencies of B1 and B2 are reversed across
592 blocks. Each block consists of 50 trials. The task requires the agent to learn that the
593 intermediate states B1 and B2, instead of the agent's own choice of A1 or A2, are actually the
594 ones that determine the reward.

595

596 The training dataset consists of trials with randomly assigned choices. This is to reflect that the
597 actual learning by human subject starts from rather random choices. During the training, only
598 the events in the rewarded trials are actually trained. This is because the goal of the training is
599 to learn the sequences more likely leading to reward. For the network to learn the event
600 contingencies across trials, the state vector L_j is not re-initialized at the beginning of each
601 trial and the events in the previous trial are included in the training procedure. The loss
602 function is calculated based only on the current rewarded trial, but the errors are back-
603 propagated to the previous trial and the network connections are updated accordingly. In the
604 testing sessions, the network connections are not further updated.

605

606 After the training, the network learns to make decisions based on the events both in the
607 current and in the previous trials. Immediately after a block change, the network performs
608 poorly because of the reversed reward contingency. Its performance then recovers gradually
609 (Fig 8a). To exclude the possibility that the network simply learns to reverse its choice every
610 50 trials, we train the network with 50-trial blocks and test it with a different block size and
611 observe similar results (Supp fig 2).

612

613 To further investigate the network's adaptive behavior, we use a factorial analysis introduced
614 by Daw and colleagues (Daw et al., 2011) to look at how different factors of the task affect
615 the choice. We sort the trials based on the intermediate outcomes and the reward outcomes
616 into four groups: common-rewarded (CR), common-unrewarded (CU), rare-rewarded (RR)
617 and rare-unrewarded (RU). Then we calculate the probability of the network repeating the
618 previous choice in the next trial. The probability is higher after the CR and RU trials than
619 after the CU and RR condition (Fig. 8b). The results indicate that the task structure
620 information is used for decisions, which was interpreted as evidence for model-based decision
621 making (Daw et al., 2011). Interestingly, even though the error signals are only propagated one
622 trial back during the training period, the history influence on the network's choice extends to
623 many trials back (Fig 8c).

624

625 Two more analyses corroborate this conclusion. We first use a logistic regression to look at
626 how a variety of factors affect network's choices (Fig 8d). The result suggests that the
627 interaction term between the intermediate outcome and the reward outcome (*Trans x Out*) is
628 the largest factor affecting the choices. We further fit the network's behavior with a mixture
629 of the model-free and model-based agent (Akam et al., 2015; Daw et al., 2011). The higher
630 weight for the model-based strategy (Table 1) indicates the network's choice behavior is closer
631 to model-based.

632

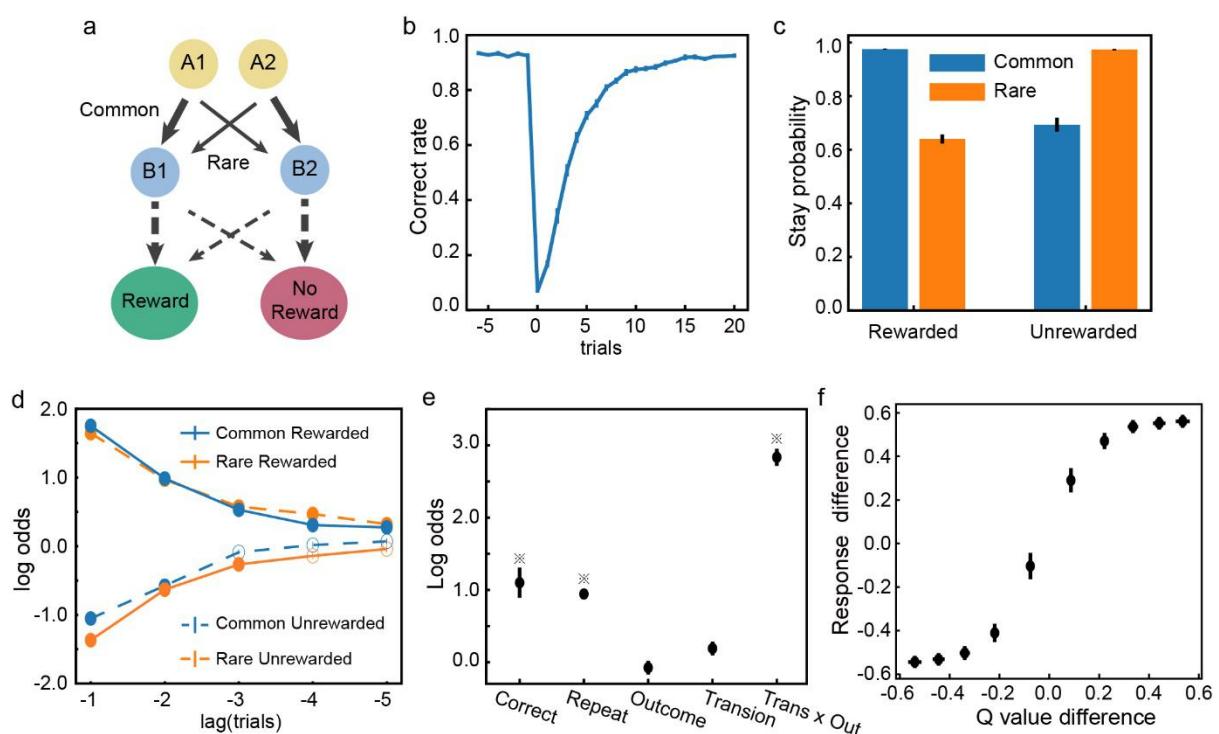
α_1	λ	α_2	γ	t
0.87 ± 0.04	0.75 ± 0.05	0.68 ± 0.05	0.47 ± 0.01	9.49 ± 1.16

633

634 Table 1. Fitted parameters. The S.E. across runs is indicated. The weight for the model-based algorithm, α_1 , is significantly larger than 0.5 ($p < 0.001$), indicating
 635 a strong contribution of the model-based behavior.

637

638 Although the network is not trained explicitly to calculate and compare the value associated
 639 with each state, it is curious whether the network activities encode the Q-value estimated with
 640 RL. Based the RL model fitting above, we estimate the choices' Q-value in each trial. When
 641 we plot the estimated Q-value difference between the two choices, it shows a clear correlation
 642 with the activity difference between the two choice output units (Figure 8e).



643

644 Figure 8. **a.** The two-step task. The thick and thin lines denote the common and rare
645 transitions, respectively. The contingencies indicated by the dashed lines are reversed
646 across blocks. **b.** The switching behavior. Trials are aligned to the block switch (trial 0).
647 The performance first drops to below the chance level but then gradually recovers. **c.** The
648 probability of repeating the previous choice. The stay probabilities of the subsequent trials
649 are higher for the CR and the RU trials than the RR and the CU trials. **d.** Trial history
650 effects. The choice in the current trial is affected by the trial types in the previous trials.
651 Solid dots indicate significant effect (Bonferroni correction, $p < 0.05$). **e.** Factors affecting
652 the choices. \ast indicates significance ($p < 0.01$). **f.** The response difference between two
653 choice output units is correlated with the difference between the estimated Q-value of the
654 two actions.

655

656

657

658 **Discussion**

659

660 Here, we present a neural network model for flexible and adaptive decision making. By
661 training essentially the same network model to perform four exemplary tasks, we show that
662 the network model not only can solve complicated decision-making problems, perform
663 Bayesian inference, and support meta-cognition, but also can use past trial history to adapt its
664 choices in a volatile environment. The units in the network exhibit signature behaviors of the
665 decision-related neurons previously described experimentally, encoding the accumulative
666 evidence and urgency signal. We also suggest a possible mechanism of the speed and accuracy
667 trade-off. The network model provides a unified solution to questions that have been
668 investigated previously with different modeling approaches.

669

670 [Drift-diffusion model and Attractor network models](#)

671

672 The probabilistic reasoning task and the post-decision wagering task have been studied with
673 the drift-diffusion model (DDM). In the DDM, the decision-making process is modeled as a
674 diffusion process in which the evidence biases the drift rate. The behavior of the DDM can be
675 accounted for by our model. This is interesting because, unlike the DDM and its variants, our

676 model does not explicitly model the decision as a competition between alternatives in which
677 evidence is accumulated. Yet by training the model to learn the statistical contingencies
678 between events, the model exhibits behavior that fits the DDM well.

679

680 Our modeling effort is a distinct approach from the previous work of attractor neural network
681 models (X.-J. Wang, 2002; Wong & Wang, 2006). Like the DDM, the attractor network
682 models treat the decision-making process as a competition. The competition in attractor
683 networks is between two groups of neurons, each receiving sensory inputs supporting its
684 corresponding option. The sensory inputs drive the network states into attractor states that
685 represent the final choice. This type of attractor network in its simplest form is obviously not
686 a general solution for more complex decision-making tasks, for example, when the decisions
687 require multiple steps as in the two step tasks, or when the contingency between the sensory
688 inputs and the reward is volatile.

689

690 Even if we limit our discussions to the kind of decision making problem that has been studied
691 with the attractor network model, the current approach provides several new insights. First,
692 our model describes a potential mechanism for the speed-accuracy tradeoff in decision
693 making distinct from the attractor models. In our model, there are a group of units that
694 contribute to reaction time causally. Through the modulation of these units, one can achieve

695 speed-accuracy tradeoff. In contrast, the speed-accuracy tradeoff in attractor network models
696 is achieved through the adjustment of the balance between excitatory and inhibitory
697 connections, which effectively modulate how fast the network may settle into an attractor state
698 (Lo, Wang, & Wang, 2015). Future experiments can be designed to distinguish these two
699 scenarios. Second, our model's decisions do not require the network to converge into attractor
700 states. This allows the model to perform the post-decision wagering task without additional
701 modules. The network's state, i.e. the units' activity level, not only reflects the decision
702 outcome, but also may be used to estimate confidence. It was found to be necessary to
703 introduce an extra layer to model such a task with a attractor network (Insabato, Pannunzi,
704 Rolls, & Deco, 2010). Finally, our model does not require its inputs to reflect evidence
705 strength, which is essential for an attractor network. Therefore, it is able to perform the
706 probability reasoning task in which the contingencies between the sensory stimuli and the
707 outcomes are arbitrary. The input flexibility allows our model to be a more general solution
708 than the attractor network models.

709

710 Bayesian Inference

711 The computation that the network carries out is essentially Bayesian inference. During
712 decision making, our network provides a representation of the probability distribution of the
713 sensory events, which are updated over the time as new inputs arrive. Unlike the probabilistic

714 population coding theory proposed previously (Ma, Beck, Latham, & Pouget, 2006), our
715 model does not depend on assumptions of the distributions of the input or the nature of the
716 contingency. The population response of the units in our model encode the entire distribution.
717 This is a more general solution than the probabilistic population coding theory. Another recent
718 model also aims to be a general solution for probabilistic computing (Orhan & Ma, 2017). In
719 comparison, our model emphasizes the temporal integration of information. The outputs of
720 our model may be directly mapped to the probability distribution, which may serve as a useful
721 top-down signal that can be sent back to the sensory cortex.

722

723 Reinforcement learning

724

725 By learning statistical relationships between sensory, action, and reward events across trials,
726 our model can exhibit an adaptive behavior, which has been mostly modeled in the field with
727 the reinforcement learning framework.

728

729 An interesting conceptual difference between our model and RL is the role of reward. The
730 reward plays a central role in RL, in which the learning is aimed at minimizing the reward
731 prediction error and maximizing the total future reward. The discrepancy between the actual
732 and the expected rewards is used to update the value of states. In contrast, our model treats the

733 reward events as part of the event sequences just as the sensory and action events. The
734 learning in our model is driven by a teaching signal that includes sensory, action, and reward
735 events. This signal could come from the dopamine system in the brain. The dopamine neurons
736 have been indicated to signal reward prediction error and used as evidence for the existence of
737 a reinforcement learning system in the brain (Schultz et al., 1997). Yet, recent findings have
738 started to reveal a wider role of dopamine neurons in signaling not just reward prediction error,
739 but also predictions of sensory stimuli or actions (Engelhard et al., 2018; Jin & Costa, 2010;
740 Wassum, Ostlund, & Maidment, 2012; Wood, Simon, Koerner, Kass, & Moghaddam, 2017).
741 These findings fit well with the current model in which the prediction errors need to be
742 calculated for not only the rewards, but also the sensory and action events. Future experiments
743 may be designed to test whether the dopamine system indeed fit this role.

744

745 Interestingly, our simulation results of the two-step task, along with several previous modeling
746 studies (J. X. Wang et al., 2018; Zhang et al., 2018), show that online updating of synaptic
747 weights is not necessary for adaptive behavior. Once the network is trained and learns the
748 contingency between the events across trials, it can generate adaptive behavior based on the
749 internal network dynamics, which encodes the past trial sequences, without further adjusting
750 its internal connection weights. The network may appear to be learning and its behavior can be

751 well explained by the reinforcement learning framework, yet no actual learning has to happen
752 at the level of neuronal connections.

753

754 [Training](#)

755

756 To train the model, it is necessary to feed the model with appropriate sequences of sensory,
757 action, and reward events. In the real brain, this means there has to be a separate mechanism
758 to generate appropriate responses at appropriate times and form these sequences in the first
759 place for the brain to learn. One solution is to start from the responses based on animals'
760 innate responses or other established stimulus-action associations relevant to the task. These
761 responses should contain a certain degree of variations that lead to different responses and
762 reward outcomes (Neuringer, 2002). For example, when we train monkeys to perform
763 complicated behavior task such as the probabilistic reasoning task, we start from the most
764 basic delayed saccade task. Monkeys have an innate tendency of making saccades toward
765 newly appearing visual stimuli at variable delays. By selectively rewarding saccades with
766 longer delays, we can train the monkeys to hold the fixation for an extended period before the
767 saccade. More components of the task can be introduced into the task gradually this way.
768 Similar strategy can be used to train our model. Starting with simple tasks that the network is
769 capable of doing, even if only occasionally, we can selectively reinforce the behavior that

770 leads to reward and use the event sequence leading to the reward for the next stage training,

771 until we reach to the final version of the task.

772

773 [Basal Ganglia](#)

774 It remains an interesting speculation which brain structures may carry out the computation as

775 our network model does. Areas such as the prefrontal cortex, the posterior parietal cortex, the

776 cerebellum, the hippocampus, and the basal ganglia receive a wide range of sensory, motor,

777 and reward inputs and are theoretically possible to carry out the computations demonstrated

778 here. Indeed, many studies of these brain structures find neurons that have similar response

779 patterns as the units in our model.

780

781 Among these brain structures, the basal ganglia are particular interesting. It has been long

782 noticed that the key feature of GRU networks, which is the gating mechanism that controls the

783 maintenance of old information and the updating with new incoming information, resembles

784 the anatomical circuitry in the basal ganglia (Frank, Loughry, & O'Reilly, 2001; O'Reilly &

785 Frank, 2006). The basal ganglia receive broad input from all over the cortex, including the

786 sensory cortices and motor cortices (Engelhard et al., 2018), which comprises all the necessary

787 information for learning sequences.

788

789 Our model may be applied to the current basal ganglia research in several distinct directions,
790 which have been largely studied separately.

791

792 First of all, our modeling results explain previous studies of the basal ganglia's role in decision
793 making. For example, Ding and Gold found caudate neurons represented accumulated
794 information in a random dot motion discrimination task (Ding & Gold, 2010, 2012). Cisek
795 and colleagues observed a group of neurons in the globus pallidus showed activity patterns
796 similar to the urgency signal (Thura & Cisek, 2017). In rodent experiments, it has been
797 reported lesions in the striatum led to deficits in evidence integration and produced a bias in
798 decision making (L. Wang, Rangarajan, Gerfen, & Krauzlis, 2018; Yartsev, Hanks, Yoon, &
799 Brody, 2018). Our model may account for the results from these studies.

800

801 Second, our model, which is trained to learn sequences, naturally explains the basal ganglia's
802 role in performing sequential actions and in procedure memory (Geddes, Li, & Jin, 2018; Jin
803 & Costa, 2010; Wood et al., 2017). It is conceivable that simply setting the input and output
804 units according to the task requirement, we may train the model to perform sequential actions.

805

806 Third, the basal ganglia have also been indicated to play a central role in habitual and goal-
807 directed behavior, both are highly relevant to our model (Graybiel, 1995; Graybiel & Grafton,

808 2015). On the one hand, by design, our model is capable of generating habitual behavior after
809 the training. The model simply follows a set sequence of events and carries out actions, which
810 may be interpreted as habitual responses, even in the case that the actions have to be
811 determined with the preceding sensory inputs in a complicated manner. On the other hand, our
812 model may also contribute to goal-directed behavior. A key component of goal-directed
813 behavior is the ability of assessing which actions may lead to a desired goal. Such assessment
814 depends on the capability of making predictions of the actions' consequences, which is exactly
815 what our model can do.

816

817 Our model may therefore help us to link the existing studies of basal ganglia from distinct
818 perspectives to form a coherent computational theory of the basal ganglia's role. Future work
819 should incorporate more details of the basal ganglia circuitry into the network structure.

820

821 [Further investigations](#)

822

823 The current results focus on the similarities between the model and the behavior and
824 physiological findings from the previous experiments. It would be interesting to further
825 explore how tweaking components and parameters of the GRU network would affect these

826 results. This knowledge should motivate the experimental investigation of the relevant neural
827 circuitry.

828

829 We also have not investigated extensively how the particular loss function we use to train the
830 network in the current study, which treats all events equally and is flat over time, can be
831 improved. One possibility is to normalize the loss function within each domain of sensory,
832 motor, and reward, and use an appropriate weight for each domain to reach a suitable balance.
833 Another possibility is to weigh loss functions differently across the time, so that events remote
834 from the reward in the sequence exerts less influence on learning. Our preliminary
835 explorations in these directions have yielded similar results (not shown), but more detailed
836 investigations may reveal a loss function that may create a network that models the brain
837 better and provides us with further insights.

838

839 Our current investigation focuses on the hidden layer, which is where the learning and the
840 decision making occurs. Conceptually, this is similar to the generative model in the predictive
841 coding framework (Friston, 2005; Rao & Ballard, 1999). The output layer in the current study
842 only serves a limited purpose for the verification of the network's performance. One may
843 extend the current work and construct a circuitry in which the model's outputs are sent to a
844 sensory module and modulate the sensory processing. The sensory module sends inputs back

845 to the model and forms a loop. This way, we can create a complete model that may be further
846 used to understand the neural circuitry in both the sensory and the motor parts of the brain.

847

848 Finally, GRU networks have been shown to be able to solve problems at much larger scales,
849 for example, the natural language processing (Chung et al., 2014). Limited by the availability
850 of the existing experimental studies, we have only tested the network with small-scale toy
851 problems. As a result, we use only a small network with 128 units. Natural language
852 processing demands the learning of complex contingency structures between elements of
853 language contained in a sequence that can be very long. Although at a much more complex
854 level, this is essentially the same learning problem as what we study with our current model.
855 By expanding our network model with more units and more layers, we should be able to model
856 much more complex behavior. It is an interesting parallelism between the studies of neural
857 network solutions for natural language processing, which develops neural networks to solve a
858 complex computational problem that the human brain does, and the brain modeling studies,
859 which uses network models to understand the computations in the brain. Our study is a
860 beginning to bridge the two fields for the benefits of both.

861 **Funding**

862

863 This work was supported by Shanghai Municipal Science and Technology Major Project

864 (Grant No.2018SHZDZX05) and by Strategic Priority Research Program of Chinese

865 Academy of Science, Grant No. XDB32070100.

866

867 **Acknowledgements**

868

869 We thank Zhongqiao Lin, Chechang Nie, Yang Xie, Wenyi Zhang for their help in all phases

870 of the study, and Shan Yu for providing comments and advice. The authors declare no

871 competing financial or nonfinancial interests.

872

873 **Methods**

874

875 [Task 1: Probabilistic reasoning task](#)

876 Network

877

878 Our model is a network that contains three layers: an input layer, a hidden layer based on gated
879 recurrent units, and an output layer (Fig 1a). The input layer contains units that carry the
880 information about the sensory, motor, and reward events. ($\mathcal{N}_B=20$).

881

882 There is a total of 14 sensory input units. 10 of them represent the 10 shapes in the task.

883 There are 3 additional units indicating the presence of the eye movement targets, including the
884 fixation point, the left target, and the right target. We also include a unit that indicates the
885 absence of any visual stimuli.

886

887 There are 4 action input units that represent the efference copies of the motor commands: 1
888 for fixation on the fixation point, 1 for saccading to and then fixating on the left target, 1 for
889 saccading to and then fixating on the right target, and 1 for saccading to other locations, which
890 is considered as a fixation break and aborts a trial.

891

892 Finally, 2 reward input units are included: one for reward and one for the absence of reward.

893

894 The output layer includes also 20 units that mirror the inputs.

895

896 We set the time step to 100ms in the model training and testing. Our analyses and conclusions

897 remain valid for smaller time steps.

898

899 Behavior Analysis

900

901 The parameters of the network are fixed during testing. We generate 5000 random shape

902 sequences of 25 shapes long as the testing data and feed them into the network model. A

903 shape sequence is stopped whenever the output units associated with the saccades are triggered.

904 If the network does not make a response before all 25 shapes have been presented or makes a

905 response at an inappropriate time, the trial is aborted (172.50 ± 3.73 or 3.45 ± 0.07 % trials in

906 each run, excluded in further analyses).

907

908 Subjective weights (Fig 2b):

909 We perform a logistic regression to assess how each kind of shape affects the choice. The

910 probability of choice is a function of the sum of leverages, Q , provided by each kind of shape:

911

912
$$\text{logLR}_n = \beta_0 + \frac{\beta_1}{4} \text{shape}_m,$$

913 (12)

914
$$\text{logLR}_n = \beta_0 + \sum_{m=1}^{10} \beta_m \text{shape}_m, \quad (13)$$

915

916 where shape_m represents the how many times shape m appears in a trial. β_0 is the bias term, β_1 to β_{10}

917 are the estimates of how much weight the network model assigns to shapes 1 to 10, and are

918 termed as the subjective weights of evidence (Yang & Shadlen, 2007). Since the regressors

919 are not independent to each other, here we use the ridge regression to minimize the variation

920 of the estimation. The hyperparameter controls the trade-off between the cross-entropy loss

921 and the L2-norm of the coefficients is selected through a ten-fold cross validation for each

922 fitting.

923

924 Shape order (Fig 2d):

925 To test how the shape order affects the network's choice, we perform a logistic regression on

926 the trials with more than six epochs:

927

928
$$\text{logLR}_n = \beta_0 + \beta_1 \text{shape}_n + \beta_2 \text{shape}_{n+1} + \dots + \beta_{n-1} \text{shape}_{n-1}, \quad (14)$$

929

930 where $\log LR_n$ is the logLR of the shape in epoch n and N is the total number of epochs in a
931 sequence. β_0 is the bias term, $\beta_1/\beta_2/\beta_3$ are the fitting coefficients of the shapes in the first
932 three epochs, β_4 the average effect for the middle epochs, and β_5/β_6 the second and the third
933 epochs to the last. The regression is done without the final shape (n -th). This is because the
934 last shape is almost always (>99% of trials) consistent with the choice. The fitting procedure is
935 similar to what we use above to estimate the subjective weight of each shape.

936

937 Unit Selectivity (Fig 3ab):

938

939 We test whether each neuron's activity in the hidden layer is modulated by the total logLR,
940 absolute value of the total logLR, the urgency, and the choice outcome with a linear regression.
941 We align the hidden unit state L to the shape onset in each epoch and use a simple linear
942 regression to characterize the unit's selectivity:

943

944
$$L_t = \beta_0 + \beta_1 \text{logLR}_{\text{tot}} + \beta_2 \text{logLR}_{\text{avg}} + \beta_3 \text{urgency} + \beta_4 \text{choice}, \quad (15)$$

945

946 where L_t is the unit's response at time t aligned to the shape onset, $\text{logLR}_{\text{tot}}$ is the total
947 logLR of the shapes that have been presented by time t , urgency represents the urgency, which is
948 quantified as the number of shape epochs, choice represents the choice of the network, which is

949 set to 1 when the left target is chosen, -1 when the right target is chosen, and 0 if fixation is
950 still maintained by the end of the epoch. The regression is performed on every time step
951 during the shape representation. We define a unit as selective to a variable if the variable
952 shows a significant effect on the unit's activity at every time step in an epoch. The significance
953 is determined by two-tailed t-tests with Bonferroni correction.

954

955 Response Variability (Fig 3d)

956

957 We calculate a unit's response variability as the standard deviation of the unit's average
958 response in each shape epoch across trials. Fig 3d plots the average response variability for all
959 units that are selective to the total logLR in all 20 simulation runs, using only the trials with
960 more than 5 shapes.

961

962 Speed-accuracy tradeoff (Fig 4d)

963

964 To test how *when* units affect the speed-accuracy tradeoff, we suppress the outputs of the
965 +*when/-when* units that are selected with variable criteria. For example, the criteria of the
966 accumulative I_{when} of the selected +*when* units are set to 10%, 20%, 30%, 40%, and 50% out
967 of the summed I_{when} of all units with positive I_{when} . At each criterium, the number of units

968 manipulated is not linearly scaled, but their total contribution to the behavior, measured by
969 their summed I_{when} , is scaled linearly.

970

971 Network analysis (Fig 4cd):

972

973 Geodesic distance and max flow are used to quantify how much information may be
974 transferred between two nodes in a graph. We use weight matrix \mathbf{L}' , which represents the
975 connection strengths between units in the hidden layer, to construct a weighted directed graph.
976 Each unit in the hidden layer corresponds to a node in the graph. Connections with the highest
977 30% of the absolute value of the connection strength are turned into the edges in the graph.
978 The results hold if we keep the highest 10% or 50%. The weight and capacity of each edge are
979 defined as the absolute value of the original connection strength.

980

981 We define the geodesic distance from node A to target B as the minimum summed inverse of
982 the weights of the edges from node A to B (Newman, 2001). The geodesic distances between
983 any pairs are calculated with Dijkstra's algorithm (Ahuja, Magnanti, & Orlin, 1993). The
984 distance between the pairs not connected is defined as infinite. To account for these pairs, we
985 compare the mean of the inverse of geodesic distance across all unit pairs in each group.

986

987 Task 2: Multi-sensory integration task

988 Input units:

989

990 There are two groups of input units representing the sensory stimuli, corresponding to the

991 visual and the vestibular inputs. Each group has 8 units with Gaussian tuning curves with

992 different preferred directions D_{pref} at -90° , -64.3° , -38.6° , -12.9° , 12.9° , 38.6° , 64.3° , and 90° ,

993 respectively. The tuning curve $f(D)$ is described as following:

994
$$f(D) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(D-D_{pref})^2}{2\sigma^2}}, \quad (16)$$

995 where D is the direction of the motion stimulus in the range between -90° and 90° and $\sigma=45^\circ$

996 is the standard deviation of the Gaussian function. The observed spike count s has a Poisson

997 distribution:

998
$$P(s) = \frac{e^{-g}}{s!} g^s, \quad (17)$$

999 where $g=300$ is a constant for the gain. Finally, to limit the response of the input units to be

1000 smaller than 1, a sigmoid-like function is used to calculate the response $r(s)$:

1001
$$r(s) = \frac{1}{1 + e^{-4(s-1)}}, \quad (18)$$

1002 The exact choice of the function is not important. Under the bi-modal condition, the responses

1003 of both the visual and the vestibular input units are calculated with the heading direction D .

1004 Under the uni-modal condition, the responses of the units of the unavailable modality are set
1005 to 0.5 without noise.

1006

1007 Bayesian inference:

1008

1009 We estimate the heading direction by calculating the discretized posterior probability of the
1010 possible motion directions;

1011 $\text{トコヒ}, \text{瞬} \angle \text{勝利} \text{管} \text{乙}). \text{乙} \text{管} \text{勝利} \text{管} \text{乙}.$ (19)

1012 The left choice target is chosen if the probability of motion direction towards left,

1013 $\text{トコヒ} \text{管} \text{乙}$, is larger than the probability of motion direction towards right, $\text{トコヒ} \text{管} \text{乙}$,

1014 Otherwise, the right target is chosen.

1015

1016 In the bi-modal condition, we integrate the information from visual and vestibular system and
1017 calculate the posterior probability $\text{トコヒ} \text{管} \text{乙} / \text{ヒ} \text{管} \text{乙}$, which is proportional to the
1018 product of the $\text{トコヒ} \text{管} \text{乙}$ and $\text{トコヒ} \text{管} \text{乙}$, given the independence between $\text{ヒ} \text{管} \text{乙}$ and
1019 $\text{ヒ} \text{管} \text{乙}$.

1020

1021 All analyses are based on 20 simulation runs.

1022

1023 Task 3: Confidence / Post-decision wagering task

1024 Input units:

1025

1026 The network has 22 input units. Ten input units are visual units that respond to motion stimuli.

1027 Five of them prefer the leftward motion and the other five prefer the rightward motion. Their

1028 internal activation, η , is linearly scaled with the coherence, and an independent Gaussian

1029 noise, $\mathcal{N}(\mu, \sigma^2)$, is added to mimic the noise during the sensory processing, where $\mu=2.5$. A

1030 sigmoid function is used as the activation function.

$$1031 \eta = \sigma \tanh\left(\frac{\mu - \text{coherence}}{\sigma}\right) + \mu, \quad (20)$$

$$1032 \eta = \sigma \tanh\left(\frac{\mu - \text{coherence}}{\sigma}\right) + \mu, \quad (21)$$

1033 where the coherence of the moving dots, coherence , is positive when the motion direction

1034 matches unit's prefer direction and negative when not.

1035

1036 The momentary evidence L_t is the response difference between the leftward-preferring

1037 units and the rightward-preferring units. The accumulated evidence E is then the sum of

1038 momentary evidence across time:

$$1039 L_t = \eta_{\text{left}} - \eta_{\text{right}}, \quad (22)$$

$$1040 E = \sum L_t, \quad (23)$$

1041 where $\epsilon_{\text{OL}}^{r/l}$ and ϵ_{RZL} are the internal activities at time point t of the i -th unit
1042 preferring the leftward and the rightward motion, respectively. $|l|$ is the duration of the motion
1043 stimulus and is randomly selected from group [1, 2, 3, 4, 6, 8] at equal probabilities. The final
1044 decision is based on the accumulated evidence \mathcal{E} . The trial sequences used to train the model is
1045 generated as follows (Kiani & Shadlen, 2009). When the sure target is not available, the left
1046 target is chosen if \mathcal{E} is larger than 0 and the right target is chosen if \mathcal{E} is smaller than 0. When
1047 the sure target is given, the left and right choice targets are selected only if $|\mathcal{E}|$ is larger than a
1048 pre-defined threshold θ . Otherwise, the sure target is selected. Furthermore, the threshold θ
1049 increases linearly over time:

1050
$$\text{Eqn 8. 5B} \quad (24)$$

1051
1052 20 independent runs are simulated for testing the reproducibility.
1053
1054 Unit activity analysis:
1055
1056 The analysis is based on choice selective units. These units are chosen based on their activities
1057 at the time point right before the choice in the trials without the sure target. The selectivity is
1058 determined by the two-tailed t-test with Bonferroni correction. The choice direction that

1059 induces a larger response is defined as a unit's preferred direction (left: $n = 46.75 \pm 1.93$, right:
1060 $n = 42.45 \pm 1.85$).

1061

1062 [Two-step Task](#)

1063 Simulation:

1064

1065 We use 20 trained session in our analysis. Each run contains 7.5×10^5 training trials (15000
1066 blocks of 50 trials) and 100 testing blocks.

1067

1068 Model fitting:

1069

1070 The analysis was previously described (Akam et al., 2015; Daw et al., 2011; Zhang et al.,
1071 2018). Briefly, the model choice is fitted to a mixed model-free and model-based algorithm.

1072 For the model-free algorithm, the value of the chosen options and observed intermediate
1073 outcomes are updated by a temporal difference algorithm with two free parameters: the
1074 learning rate α_1 and the eligibility λ , representing the proportion of the reward prediction error
1075 that is attributed to the first-stage chosen options A1 and A2. In the model-based algorithm,
1076 the value of the observed intermediate outcome is also learned with a temporal difference
1077 algorithm. The value of each option is the sum of the products of the transition probabilities

1078 and the values of the intermediate outcomes. An extra parameter, learning rate α_2 , is used for
1079 the update of the value of the intermediate outcomes in the model-based algorithm. The
1080 overall value of each option is the weighted average of its value calculated with the model-free
1081 and model-based algorithm ($w_{model-free} + w_{model-based}=1$). Finally, the probability of choosing each
1082 option is a softmax function of the values of the two options:

1083

1084
$$P_{+z} = \frac{e^{V_{model-free}(+z) + \beta V_{model-based}(+z)}}{e^{V_{model-free}(+z) + \beta V_{model-based}(+z)} + e^{V_{model-free}(+z^*) + \beta V_{model-based}(+z^*)}}, \quad (25)$$

1085

1086 where inverse temperature parameter, β controls the randomness of the choice. P_{+z} is
1087 set to 1 if action $+z$ is chosen in the previous trial and the parameter p captures the tendency
1088 of repeating the previous trial. V_z is the value of option $+z$. Together, there are six free
1089 parameters and a maximum likelihood estimation algorithm is used for fitting.

1090

1091 Logistic regression:

1092

1093 The analysis was described previously (Akam et al., 2015; Zhang et al., 2018). Briefly, five
1094 potential factors are tested with a logistic regression, which are *Correct*—a tendency to choose
1095 the choice with higher reward probability; *Repeat*—a tendency to repeat the choice no matter

1096 what the reward outcome is; *Outcome*—a tendency to repeat the rewarded choice in the
1097 previous trial; *Transition*—a tendency to repeat the choice when it leads to common
1098 intermediate outcome and switch when rare intermediate outcome is represented; *Trans x*
1099 *Out*—a tendency to repeat the same choice when the previous trial is CR or RU, and to switch
1100 the choice if the previous trial is CU or RR.
1101
1102
1103

1104 **References**

- 1105 Ahuja, R. K., Magnanti, T. L., & Orlin, J. B. (1993). *Network Flows: Theory, Algorithms, and*
1106 *Applications* | Pearson. Prentice Hall.
- 1107 Akam, T., Costa, R., & Dayan, P. (2015). Simple Plans or Sophisticated Habits? State,
1108 Transition and Learning Interactions in the Two-Step Task. *PLOS Computational Biology*,
1109 11(12), e1004648-25.
- 1110 Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., &
1111 Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder-Decoder for
1112 Statistical Machine Translation.
- 1113 Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical Evaluation of Gated
1114 Recurrent Neural Networks on Sequence Modeling.
- 1115 Churchland, A. K., Kiani, R., Chaudhuri, R., Wang, X.-J., Pouget, A., & Shadlen, M. N.
1116 (2011). Variance as a Signature of Neural Computations during Decision Making.
1117 *Neuron*, 69(4), 818–831.
- 1118 Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-Based
1119 Influences on Humans' Choices and Striatal Prediction Errors. *Neuron*, 69(6), 1204–1215.
- 1120 Dayan, P., & Niv, Y. (2008). Reinforcement learning: The Good, The Bad and The Ugly.
1121 *Current Opinion in Neurobiology*, 18(2), 185–196.
- 1122 <https://doi.org/10.1016/J.CONB.2008.08.003>

- 1123 Ding, L., & Gold, J. I. (2010). Caudate encodes multiple computations for perceptual
- 1124 decisions. *The Journal of Neuroscience : The Official Journal of the Society for*
- 1125 *Neuroscience*, 30(47), 15747–15759.
- 1126 Ding, L., & Gold, J. I. (2012). Separate, causal roles of the caudate in saccadic choice and
- 1127 execution in a perceptual decision task. *Neuron*, 75(5), 865–874.
- 1128 Engelhard, B., Finkelstein, J., Cox, J., Fleming, W., Jang, H. J., Ornelas, S., ... Witten, I.
- 1129 (2018). Specialized and spatially organized coding of sensory, motor, and cognitive
- 1130 variables in midbrain dopamine neurons. *BioRxiv*, 456194.
- 1131 <https://doi.org/10.1101/456194>
- 1132 Frank, M. J., Loughry, B., & O'Reilly, R. C. (2001). Interactions between frontal cortex and
- 1133 basal ganglia in working memory: a computational model. *Cognitive, Affective &*
- 1134 *Behavioral Neuroscience*, 1(2), 137–160.
- 1135 Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the Royal*
- 1136 *Society B: Biological Sciences*, 360(1456), 815–836.
- 1137 <https://doi.org/10.1098/rstb.2005.1622>
- 1138 Geddes, C. E., Li, H., & Jin, X. (2018). Optogenetic Editing Reveals the Hierarchical
- 1139 Organization of Learned Action Sequences. *Cell*, 174(1), 32-43.e15.
- 1140 <https://doi.org/10.1016/j.cell.2018.06.012>
- 1141 Graybiel, A. M. (1995). Building action repertoires: memory and learning functions of the

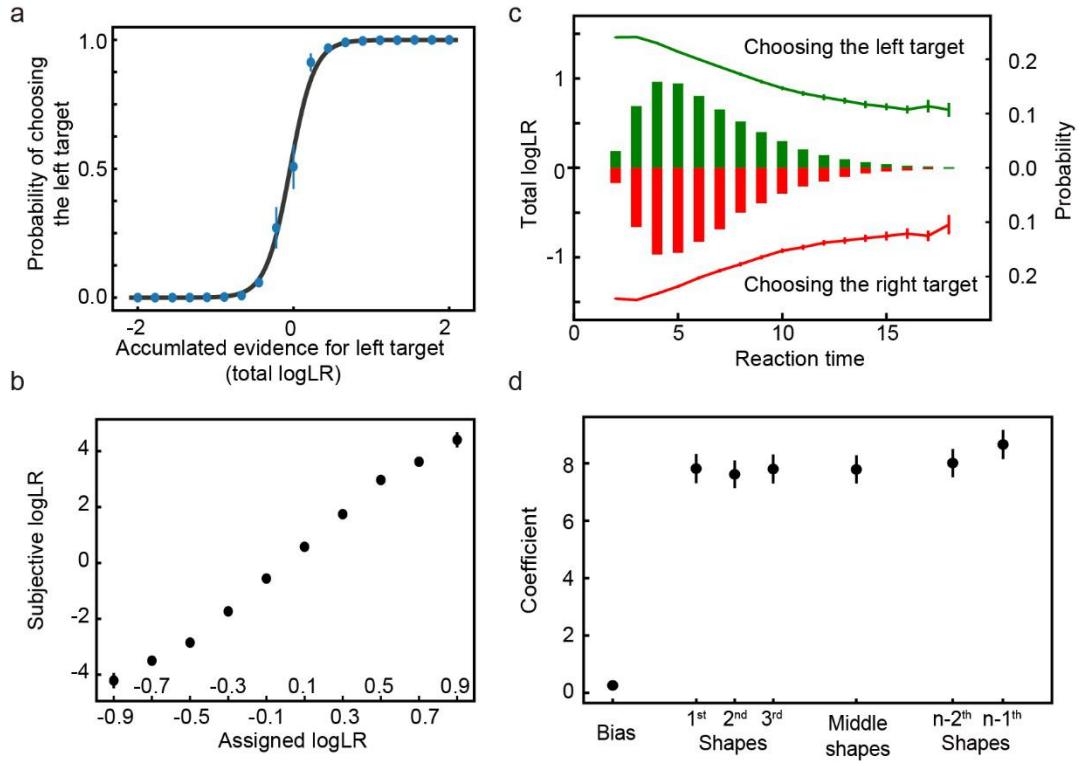
- 1142 basal ganglia. *Current Opinion in Neurobiology*, 5(6), 733–741.
- 1143 [https://doi.org/10.1016/0959-4388\(95\)80100-6](https://doi.org/10.1016/0959-4388(95)80100-6)
- 1144 Graybiel, A. M., & Grafton, S. T. (2015). The striatum: where skills and habits meet. *Cold*
- 1145 *Spring Harbor Perspectives in Biology*, 7(8), a021691.
- 1146 <https://doi.org/10.1101/cshperspect.a021691>
- 1147 Greff, K., Srivastava, R. K., & Koutník, J. (2016). LSTM: A search space odyssey. *IEEE*
- 1148 *Transactions on ...*, 1–11.
- 1149 Gu, Y., Angelaki, D. E., & DeAngelis, G. C. (2008). Neural correlates of multisensory cue
- 1150 integration in macaque MSTd. *Nature Neuroscience*, 11(10), 1201–1210.
- 1151 <https://doi.org/10.1038/nn.2191>
- 1152 Hanks, T., Kiani, R., & Shadlen, M. N. (2014). A neural mechanism of speed-accuracy
- 1153 tradeoff in macaque area LIP. *ELife*, 3.
- 1154 Heitz, R. P., & Schall, J. D. (2012). Neural mechanisms of speed-accuracy tradeoff. *Neuron*.
- 1155 Insabato, A., Pannunzi, M., Rolls, E. T., & Deco, G. (2010). Confidence-Related Decision
- 1156 Making. *Journal of Neurophysiology*, 104(1), 539–547.
- 1157 <https://doi.org/10.1152/jn.01068.2009>
- 1158 Jin, X., & Costa, R. M. (2010). Start/stop signals emerge in nigrostriatal circuits during
- 1159 sequence learning. *Nature*, 466(7305), 457–462. <https://doi.org/10.1038/nature09263>
- 1160 Kiani, R., & Shadlen, M. N. (2009). Representation of Confidence Associated with a Decision

- 1161 by Neurons in the Parietal Cortex. *Science*, 324(5928), 759–764.
- 1162 <https://doi.org/10.1126/science.1169405>
- 1163 Kim, J.-N., & Shadlen, M. N. (1999). Neural correlates of a decision in the dorsolateral
1164 prefrontal cortex of the macaque. *Nature Neuroscience*, 2(2), 176–185.
- 1165 <https://doi.org/10.1038/5739>
- 1166 Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization.
- 1167 Kira, S., Yang, T., & Shadlen, M. N. N. (2015). A neural implementation of Wald's sequential
1168 probability ratio test. *Neuron*, 85(4), 861–873.
- 1169 <https://doi.org/10.1016/j.neuron.2015.01.007>
- 1170 Lo, C.-C., Wang, C.-T., & Wang, X.-J. (2015). Speed-accuracy tradeoff by a control signal
1171 with balanced excitation and inhibition. *Journal of Neurophysiology*, 114(1), 650–661.
- 1172 <https://doi.org/10.1152/jn.00845.2013>
- 1173 Ma, W. J., Beck, J. M., Latham, P. E., & Pouget, A. (2006). Bayesian inference with
1174 probabilistic population codes. *Nature Neuroscience*, 9(11), 1432–1438.
- 1175 Neuringer, A. (2002). Operant variability: Evidence, functions, and theory. *Psychonomic
1176 Bulletin & Review*, 9(4), 672–705. <https://doi.org/10.3758/BF03196324>
- 1177 Newman, M. E. J. (2001). Scientific collaboration networks . II . Shortest paths , weighted
1178 networks , and centrality, 64, 1–7. <https://doi.org/10.1103/PhysRevE.64.016132>
- 1179 O'Reilly, R. C., & Frank, M. J. (2006). Making working memory work: A computational

- 1180 model of learning in the prefrontal cortex and basal ganglia. *Neural Computation*, 18(2),
- 1181 283–328. <https://doi.org/10.1162/089976606775093909>
- 1182 Orhan, A. E., & Ma, W. J. (2017). Efficient probabilistic inference in generic neural networks
- 1183 trained with non-probabilistic feedback. *Nature Communications*, 1–14.
- 1184 Rao, R. P. N., & Ballard, D. (1999). Predictive coding in the visual cortex: a functional
- 1185 interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1),
- 1186 79–87. <https://doi.org/10.1038/4580>
- 1187 Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85(2), 59–108.
- 1188 <https://doi.org/10.1037/0033-295X.85.2.59>
- 1189 Roitman, J. D., & Shadlen, M. N. (2002). Response of Neurons in the Lateral Intraparietal
- 1190 Area during a Combined Visual Discrimination Reaction Time Task. *Journal of*
- 1191 *Neuroscience*, 22(21), 9475–9489.
- 1192 Schultz, W., Dayan, P., & Montague, P. R. (1997). Neural Substrate of Prediction and. *Science*,
- 1193 275(5306), 1593–1599.
- 1194 Stone, M. (1960). Models for choice-reaction time. *Psychometrika*, 25(3), 251–260.
- 1195 <https://doi.org/10.1007/BF02289729>
- 1196 Sutton, R. S., & Barto, A. G. (2012). Reinforcement Learning: An Introduction, 1–334.
- 1197 Thura, D., & Cisek, P. (2017). The Basal Ganglia Do Not Select Reach Targets but Control
- 1198 the Urgency of Commitment. *Neuron*, 95(5), 1160-1170.e5.

- 1199 <https://doi.org/10.1016/j.neuron.2017.07.039>
- 1200 Wald, A., & Wolfowitz, J. (1948). Optimum Character of the Sequential Probability Ratio
1201 Test. *The Annals of Mathematical Statistics*, 19(3), 326–339.
- 1202 <https://doi.org/10.1214/aoms/1177730197>
- 1203 Wang, J. X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z., ...
- 1204 Botvinick, M. (2018). Prefrontal cortex as a meta-reinforcement learning system. *Nature
1205 Neuroscience*, 21(6), 860–868. <https://doi.org/10.1038/s41593-018-0147-8>
- 1206 Wang, L., Rangarajan, K. V., Gerfen, C. R., & Krauzlis, R. J. (2018). Activation of Striatal
1207 Neurons Causes a Perceptual Decision Bias during Visual Change Detection in Mice.
Neuron, 97(6), 1369-1381.e5. <https://doi.org/10.1016/J.NEURON.2018.01.049>
- 1209 Wang, X.-J. (2001). Synaptic reverberation underlying mnemonic persistent activity. *Trends in
1210 Neurosciences*, 24(8), 455–463.
- 1211 Wang, X.-J. (2002). Probabilistic decision making by slow reverberation in cortical circuits.
Neuron, 36(5), 955–968.
- 1213 Wassum, K. M., Ostlund, S. B., & Maidment, N. T. (2012). Phasic Mesolimbic Dopamine
1214 Signaling Precedes and Predicts Performance of a Self-Initiated Action Sequence Task.
Biological Psychiatry, 71(10), 846–854.
- 1216 <https://doi.org/10.1016/J.BIOPSYCH.2011.12.019>
- 1217 Wong, K.-F., & Wang, X.-J. (2006). A Recurrent Network Mechanism of Time Integration in

- 1218 Perceptual Decisions. *Journal of Neuroscience*, 26(4), 1314–1328.
- 1219 <https://doi.org/10.1523/JNEUROSCI.3733-05.2006>
- 1220 Wood, J., Simon, N. W., Koerner, F. S., Kass, R. E., & Moghaddam, B. (2017). Networks of
1221 VTA Neurons Encode Real-Time Information about Uncertain Numbers of Actions
- 1222 Executed to Earn a Reward. *Frontiers in Behavioral Neuroscience*, 11, 140.
- 1223 <https://doi.org/10.3389/fnbeh.2017.00140>
- 1224 Yang, T., & Shadlen, M. N. (2007). Probabilistic reasoning by neurons. *Nature*, 447(7148),
1225 1075–1080. <https://doi.org/10.1038/nature05852>
- 1226 Yartsev, M. M., Hanks, T. D., Yoon, A. M., & Brody, C. D. (2018). Causal contribution and
1227 dynamical encoding in the striatum during evidence accumulation. *eLife*, 7, 1–24.
- 1228 <https://doi.org/10.7554/eLife.34929>
- 1229 Zhang, Z., Cheng, Z., Lin, Z., Nie, C., & Yang, T. (2018). A neural network model for the
1230 orbitofrontal cortex and task space acquisition during reinforcement learning. *PLOS*
1231 *Computational Biology*, 14(1), e1005925. <https://doi.org/10.1371/journal.pcbi.1005925>
- 1232
- 1233



1234

1235 Supplementary Figure 1. The performance of the network trained with a dataset that

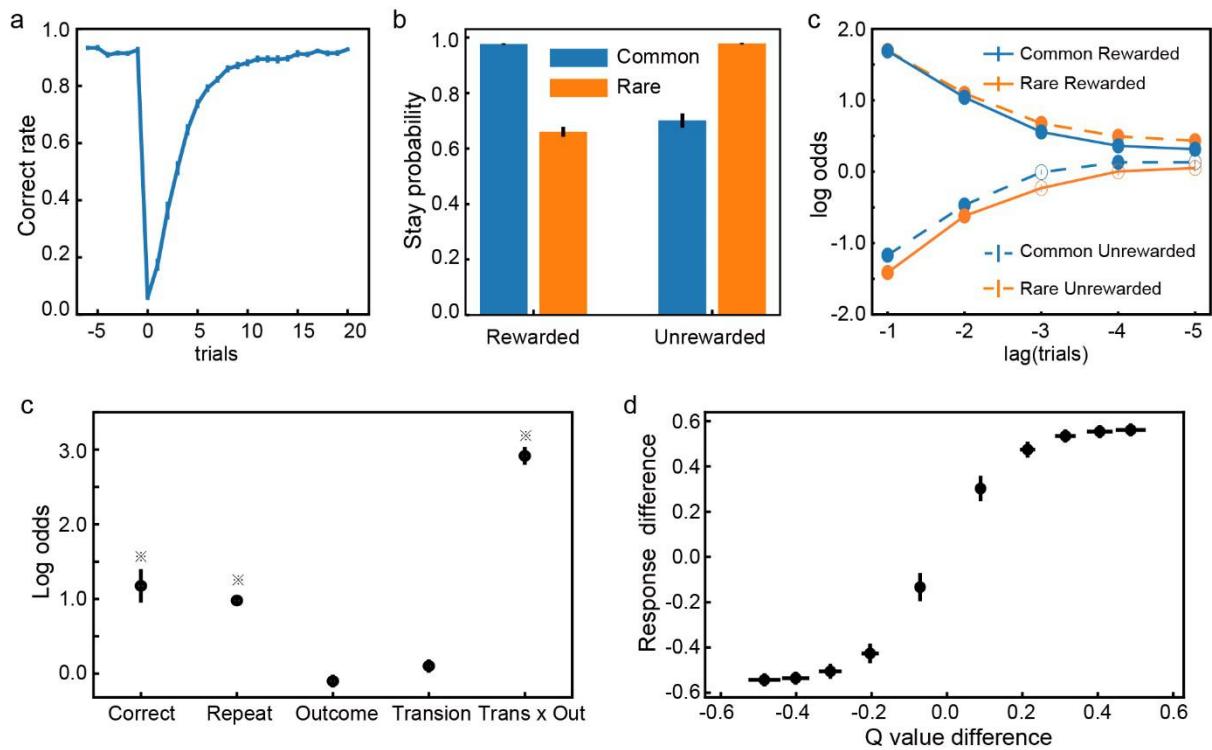
1236 contains only 1000 unique sequences. Same format as in Figure 2.

1237

1238

1239

1240



1241

1242 Supplementary Figure 2. The network's performance when the testing dataset's block size
1243 is set to 70 trials. Same format as in Figure 8.

1244