

1
2
3
4
5 **Context-dependent persistency**

6 **as a coding mechanism for robust and**
7 **widely distributed value coding**

8
9
10
11
12 Ryoma Hattori¹ and Takaki Komiyama^{1,2}

13
14
15
16 ¹Neurobiology Section, Center for Neural Circuits and Behavior, Department of
17 Neurosciences, and Halıcıoğlu Data Science Institute, University of California San
18 Diego, La Jolla, CA 90093, USA

19 ²Lead Contact

20 *Correspondence: rhattori0204@gmail.com (R.H.), tkomiyama@ucsd.edu (T.K.)

21 **SUMMARY**

22 Task-related information is widely distributed across the brain with different coding properties,
23 such as persistency. We found in mice that coding persistency of action history and value was
24 variable across areas, learning phases, and task context, with the highest persistency in the
25 retrosplenial cortex of expert mice performing value-based decisions where history needs to be
26 maintained across trials. Persistent coding also emerged in artificial networks trained to perform
27 mouse-like reinforcement learning. Persistency allows temporally untangled value
28 representations in neuronal manifolds where population activity exhibits cyclic trajectories that
29 transition along the value axis after action outcomes, collectively forming cylindrical dynamics.
30 Simulations indicated that untangled persistency facilitates robust value retrieval by downstream
31 networks. Even leakage of persistently maintained value through non-specific connectivity could
32 contribute to the brain-wide distributed value coding with different levels of persistency. These
33 results reveal that context-dependent untangled persistency facilitates reliable signal coding and
34 its distribution across the brain.

35

36 **INTRODUCTION**

37 The parallel distributed processing (PDP) theory (McClelland et al., 1986; Rogers and
38 McClelland, 2014; Rumelhart et al., 1986) highlights computational advantages of distributed
39 information coding in neural networks and has had a profound impact on our understanding of
40 cognition and deep learning. Growing evidence revealed that information coding in the brain is
41 highly distributed across neurons and distinct brain areas (Allen et al., 2019; Hattori et al., 2019;
42 Koay et al., 2020; Musall et al., 2019; Steinmetz et al., 2019; Stringer et al., 2019). Even
43 neurons in the primary sensory cortex, which were classically thought to process only sensory
44 information of a single modality, have been found to encode diverse information such as other
45 sensory modalities (Hattori and Hensch, 2017; Hattori et al., 2017; Iurilli et al., 2012),
46 spontaneous movements (Musall et al., 2019; Stringer et al., 2019), actions (Hattori et al., 2019;
47 Koay et al., 2020; Steinmetz et al., 2019), reward (Hattori et al., 2019; Koay et al., 2020), event
48 history (Hattori et al., 2019; Koay et al., 2020), and value (Hattori et al., 2019; Serences, 2008).
49 Although these signals are widely distributed, activity perturbations of a brain area typically
50 affect only a subset of behavioral outputs that are associated with the information encoded in
51 the area. These results suggest that, although information coding is highly distributed, not all of
52 the information represented in neural activity may be used in each area.

53 A clue to understand the utility of encoded information may lie in the temporal dynamics
54 of the information coding. In working memory tasks where information is maintained for several
55 seconds in a trial, information can be maintained as either persistent neural activity or
56 sequential transient activity across a neural population that tiles the memory period (Cavanagh
57 et al., 2018; Fuster and Alexander, 1971; Masse et al., 2019; Miller et al., 1996; Murray et al.,
58 2017; Orhan and Ma, 2019; Romo et al., 1999; Zhu et al., 2020). Recently, it was shown that
59 certain brain areas in mice such as the retrosplenial cortex (RSC) (Hattori et al., 2019) and the
60 medial prefrontal cortex (Bari et al., 2019) encode action values with exceptional persistency
61 during history-dependent value-based decision making tasks where values need to be stably
62 maintained across trials. Inactivation of either area impaired the ability to use the action value
63 for their decision making. These results suggest that persistent value coding is critical for
64 animals to exploit value for decision making when the value needs to be maintained for
65 extended periods of time. Similar persistent coding is prevalent across the brain and species,
66 ranging from coding of motor planning (Guo et al., 2017; Inagaki et al., 2019; Li et al., 2016),
67 internal states (Allen et al., 2019; Marques et al., 2020) to emotions (Jung et al., 2020; Kennedy
68 et al., 2020), yet the computational advantages of persistent coding has not been fully
69 established quantitatively.

70 Here we investigated the neural dynamics of action history and value coding in 6 areas
71 of the mouse cortex and artificial recurrent neural network (RNN) agents to understand the
72 computational advantages of persistent coding and its impact on distributed coding.

73

74 **RESULTS**

75 **Learning- and context-dependence of coding persistency across cortical areas**

76 We first used the neural activity data recorded in mice performing decision making based on
77 history-dependent action value we reported previously (Hattori et al., 2019). Each trial consisted
78 of a ready period, an answer period, and an inter-trial-interval (ITI). The duration of each period
79 was variable from trial to trial, making the task more naturalistic than a fixed temporal sequence
80 (Figure 1A). During the ready period (LED cue), mice needed to withhold licking to enter the
81 answer period. This ensured that the neural activity during the ready period was free of licking-
82 related motor activity. Mice were allowed to freely choose either left or right lickport after a go
83 cue tone. Different reward probabilities were assigned to the 2 lickports, and the reward
84 probabilities changed every 60-80 trials without cue. Therefore, mice were encouraged to
85 dynamically estimate the underlying reward probabilities of the 2 options on a trial-by-trial basis

86 by forming subjective action values based on their recent choice outcome history using
87 reinforcement learning (RL) (Sutton and Barto, 2018). The action values need to be stably
88 maintained within each trial and updated after each trial based on the action and its outcome.
89 Neural activity was collected with *in vivo* 2-photon calcium imaging from transgenic mice that
90 express GCaMP6s (Chen et al., 2013) in excitatory neurons (Wekselblatt et al., 2016) (Figure
91 1B), and the calcium signals were converted to estimated spike rates by non-negative
92 deconvolution (Friedrich et al., 2017; Pachitariu et al., 2018). The recording data were from 6
93 cortical areas including 2 association (RSC: retrosplenial; PPC: posterior parietal), 2 premotor
94 (pM2: posterior secondary motor; ALM: anterior-lateral motor), and 2 primary sensory (S1:
95 primary somatosensory; V1: primary visual) cortex. We estimated the 2 action values on each
96 trial (Q_L and Q_R) by fitting a RL model to the choices of mice, and we focused our analyses on
97 the neural coding of the policy value ($\Delta Q = Q_L - Q_R$: the value difference between the 2 actions)
98 on which animals rely their decision making.

99 Regression analysis of the activity of individual neurons at different time bins within the
100 ready period identified significant fractions of neurons that encode ΔQ in all 6 areas, with the
101 highest fraction in RSC (Figure 1C). ΔQ coding in these neurons was independent of upcoming
102 choice directions (Figure S1), and reliably updated at single-trial resolution (Figure S2),
103 indicating that these neurons faithfully encoded ΔQ on a trial-by-trial basis across all 6 areas.
104 Despite the widespread ΔQ coding, the temporal stability of ΔQ coding within the ready period
105 differed across areas. Only in RSC, the ΔQ -coding neurons identified at different time bins
106 reliably encoded ΔQ throughout the trial and across trials, while the encoding was temporally
107 unstable in the other 5 areas (Figure 1D and S1H). This was because the way individual
108 neurons encoded ΔQ across time differed across areas (Figure 1E and S1I). We quantified the
109 temporal stability of ΔQ coding by defining the persistency index which reflects the coding
110 persistency relative to the chance level (Methods). The analysis revealed RSC as the area with
111 the highest ΔQ coding persistency (Figure 1F).

112 We next examined whether the coding persistency is a fixed property of individual areas
113 or changes with learning. We analyzed the population activity from RSC, PPC, pM2 and ALM
114 during early stages of training (< 1 week from training start, Figure 1G). We compared their
115 value coding persistency between early and expert sessions. We found that the ΔQ coding
116 persistency significantly increases in RSC, PPC and pM2 during training (Figure 1H-I),
117 indicating that coding persistency can change during task learning.

118 The coding persistency may have increased during learning because the value-based
119 decision task requires stable value maintenance for an extended period of time across trials.
120 Therefore, we tested whether coding persistency differs in another task that does not require
121 long maintenance of value. We trained 9 mice in the alternate choice task in which a reward
122 was given when mice made a choice that was the opposite to the previous action (Figure 2A).
123 Thus, the correct action depended only on the immediately preceding trial, in contrast to the
124 value task in which history from multiple past trials was informative. All other task conditions
125 were identical between the 2 tasks. *camk2-tTA::tetO-GCaMP6s* transgenic mice were trained in
126 the alternate choice task for at least 2 weeks to achieve a plateau-level performance (~80%
127 correct) (Figure 2B). We then performed 2-photon calcium imaging of 8,524 RSC cells, 3,186
128 PPC cells, 7,915 pM2 cells and 4,911 ALM cells (RSC: 14 populations, 608.9 ± 18.1 cells, PPC:
129 7 populations, 455.1 ± 25.1 cells, pM2: 14 populations, 565.4 ± 34.6 cells, ALM: 10 populations,
130 491.1 ± 36.8 cells, mean \pm s.e.m per population). The coding persistency of action history in the
131 alternate choice task was significantly weaker than in value-based task for the 4 imaged areas
132 (Figure 2C, D). These results indicate that the coding persistency in the cortex is context-
133 dependent.

134

135 **Persistent value coding in RSC forms cylindrical dynamics**

136 In the value-based task, ΔQ coding in RSC is temporally stable within each trial. However, this
137 does not necessarily mean that RSC population activity is static during these periods. In fact,
138 individual neurons in RSC showed heterogeneous and rather dynamic activity patterns (Figure
139 1B). To investigate how the coding of different information temporally interacts, we sought to
140 decompose population activity into the demixed neural subspaces where different task-related
141 signals are separated into distinct dimensions. Specifically, we sought to define 3 demixed axes
142 each encoding ΔQ , Q_{ch} (value of selected action, e.g. Q_L on left choice trial), or ΣQ (sum of 2
143 values), and the remaining Q -free subspace that retains all the activity variance that is not
144 explained by the 3 Q -related axes. A previous study reported demixed principal component
145 analysis (dPCA) (Kobak et al., 2016) as a method to decompose population activity into
146 demixed target-dependent and independent dimensions. However, dPCA is only designed to
147 identify dimensions for discrete variables and cannot be applied for continuous variables such
148 as Q -related signals. In addition, dPCA splits each targeted signal into multiple linear axes,
149 which makes the signal interpretation difficult. To overcome these limitations, we developed a
150 novel dimensionality reduction method that is more generally applicable, which we term

151 demixed subspace principal component analysis (dsPCA) (Figure 3A). dsPCA identifies
152 demixed dimensions for targeted signals and dimensions for target-independent activity,
153 similarly to dPCA. However, unlike dPCA, it groups each of the target signals along a single
154 linear coding dimension and can identify such dimensions for both discrete and continuous
155 target variables. The first step of dsPCA identifies the best demixed linear axes for the target
156 variables using a regression-based approach, similarly to (Mante et al., 2013). This step
157 involves fitting a multiple linear regression model of the form $x(trial) = \beta_A A(trial) +$
158 $\beta_B B(trial) + \beta_C C(trial) + \beta_0$ to the activity of individual neurons for the targeted variables, A, B
159 and C. The regression coefficients, β_A , β_B and β_C are the partial derivatives of the neural activity
160 by each target variable, and the vectors that consist of the coefficients from all neurons are the
161 linearly demixed coding directions of the neural population for the 3 targeted variables. We
162 defined the targeted coding axes as the unit vectors of these coding directions. By definition,
163 these demixed coding vectors capture all linear information of targeted variables in a population.
164 Next, dsPCA identifies the remaining target-free subspace that is orthogonal to these targeted
165 axes and captures all the remaining activity variance. The target-free orthogonal subspace is
166 identified by performing full QR decomposition of the matrix with the coding axis vectors. Then
167 the axes of the target-free subspace are further realigned based on the principal components of
168 the activity within the target-free subspace to define axes that contain large fractions of
169 remaining variance. (Figure 3B). Therefore, dsPCA can be viewed as a general extension of
170 PCA by combining the regression-based supervised target axis identifications and the PCA-
171 based unsupervised dimensionality reduction of the target-free population dynamics.

172 We evaluated the demixing performance of dsPCA using noisy simulated neural
173 populations (200 neurons / population with Gaussian noise) where graded signals A, B and C
174 are linearly encoded in 20% of the neurons. Each target signal was uniquely encoded only
175 along the single, target axis (Figure 3C-D), and linear decoders failed to decode any A, B and C
176 signals in the remaining target-free subspace (Figure 3E). We next applied dsPCA on the
177 cortical population activity time-averaged over the ready period to identify demixed coding axes
178 for ΔQ , Q_{ch} , and ΣQ , and the remaining, Q -free subspace (Figure 3F). For all 6 areas, most of
179 the targeted information was confined to each of the coding axes, and the remaining subspace
180 completely lacked any of the targeted information even though this subspace contained the
181 highest activity variance (Figure 3G-I, and S3). Although we detected some Q_{ch} signal along the
182 ΣQ axis (Figures 3H and S3B), this is expected because Q_{ch} is a component of ΣQ ($\Sigma Q = Q_{ch} +$
183 unchosen Q). However, note that ΣQ signal is not detectable along the Q_{ch} axis, indicating that

184 the demixing of activity variance worked correctly. Thus, dsPCA successfully identified demixed
185 coding axes for Q-related variables and the remaining Q-free subspace.

186 With dsPCA, we examined how ΔQ coding temporally interacts with other dynamics. The
187 activity dynamics around the choices (between ± 4 sec from the choice) was visualized in the
188 neuronal manifold consisting of the ΔQ coding axis and the other value-related axes (Figure 3J),
189 or the manifold consisting of the ΔQ coding axis and 2 largest temporal activity variance axes
190 within the Q-free subspace (Figures 3K). We found in both manifolds that activity trajectories in
191 RSC from trials with different ΔQ values do not cross with each other across time. In the
192 manifold with the largest temporal dynamics (Figures 3K, S4 and S5), RSC population remained
193 in the initial positions linearly segregated along ΔQ axis according to ΔQ of the trial ('Pre-choice'
194 in the figures). Around the go cue time, the RSC population diverged from these initial positions
195 and drew rotational dynamics. After a choice, the population returned towards the initial
196 positions following a circular geometry. The return geometry was warped along ΔQ axis,
197 reflecting the reward prediction error (RPE) on each trial depending on the choice and its
198 outcome, which updates the ΔQ representation in the population (Figure 3L-M). The RPE-
199 dependent, bidirectional transition of the activity state ensures that the neural population closely
200 represents and updates the ΔQ coding online in each trial. In contrast, the dynamics in S1 and
201 V1 were highly tangled over time, and similar ΔQ values could accompany different activity
202 states at different time. Therefore, although ΔQ coding is widely distributed across the cortex,
203 the different levels of persistency confer different levels of tangling in ΔQ coding (Figure 3N).
204 The exceptionally high ΔQ coding persistency in RSC allows a temporally untangled value
205 representation with the within-trial cyclic dynamics that transitions along the value axis to reflect
206 value updates. These dynamics across trials collectively form cylindrical dynamics during task
207 performance.

208

209 **Untangled, persistent value coding emerges in the RNN trained with the mouse RL
210 strategy**

211 The persistent and untangled ΔQ coding in RSC, together with our previous observation that
212 RSC inactivation impairs value-based decision (Hattori et al., 2019), raises the possibility that
213 persistent value coding is advantageous in the task. We investigated this possibility by training
214 artificial RNN agents to perform RL in the same task and subsequently examining the ΔQ
215 coding scheme in the trained network. The training of RNNs was done without constraining the

216 activity dynamics. We reasoned that, if persistent coding is advantageous, trained RNN agents
217 may use persistent coding to perform the task.

218 First, we trained RNNs to perform RL optimally by teaching them the ideal choices of
219 each trial based on the reward assignment rule. In this task, once a reward is assigned to a
220 choice, the reward remains assigned until the choice is selected. As a result, the actual reward
221 probability of a choice cumulatively increases if the choice is not selected in the recent trials.
222 Therefore, an optimal choice would depend on the current reward assignment probabilities,
223 which are unknown to mice and RNN agents, and past choice history. By using the optimal
224 choices as the teacher, we trained synaptic weights of RNNs such that the RNNs use only
225 history of choice and reward to make near-optimal decisions (Figures 4A and 4B). The durations
226 between decisions were made variable, similarly to the task structure in mice. The RNNs
227 receive action outcome information only at the time step after choice and need to maintain the
228 information through recurrent connectivity across time steps and trials. These optimally trained
229 networks (“optimal RNN agents”) achieved higher reward rate than expert mice (Figure 4E).
230 Furthermore, the choice patterns of optimal RNN agents diverged from the RL model that has
231 been optimized to describe the behavior of expert mice (Figure 4E), indicating that the optimal
232 RNN acquired a RL strategy that is distinct from mice. Accordingly, a regression analysis
233 showed that the dependence of optimal RNN agents on choice and reward history differed from
234 that of expert mice (Figure 4F).

235 To obtain a network model that better mimics the mouse strategy, we trained RNNs to
236 imitate expert mouse behaviors using behavioral cloning, a form of imitation learning (Osa et al.,
237 2018). We used 50,472 decision making trials of expert mice as the teaching labels to train the
238 synaptic weights of the RNN. The goal of this training was for the RNN to make the same
239 decisions as expert mice with its recurrent activity dynamics based on the same history of
240 choice and outcome in the past trials (Figure 4C). The trained RNNs (“mouse-like RNN agents”)
241 performed RL using their recurrent activity (Figure 4D), and the reward rate and the RL model
242 accuracy were equivalent to those of expert mice (Figure 4E). Furthermore, the mouse-like RNN
243 agents used history from previous trials for its decisions in a similar way as expert mice (Figure
244 4F). Therefore, the RL strategy of expert mice was successfully transferred to the synaptic
245 weights of the trained RNN agents, and the trained RNNs could implement mouse-like RL using
246 its recurrent activity dynamics without updating synaptic weights from trial to trial.

247 We then examined how the mouse-like RNN agents encoded ΔQ . We found that RSC-
248 like persistent ΔQ coding emerged in their recurrent activity (Figures 4G). This observation is

249 significant as the training procedure did not impose *a priori* constraints on the coding scheme of
250 the RNN. We also examined how the population activity dynamics evolved during training. We
251 had RNN agents at 3 stages of training run the task (before training, intermediate (after 1 epoch
252 of training), and fully trained) and analyzed their recurrent activity during the task performance.
253 dsPCA revealed that untrained networks with random connectivity exhibit highly tangled ΔQ
254 coding, while training gradually shaped the networks to form stacked circular dynamics (Figure
255 5A). Unlike RSC that formed cylindrical dynamics (Figure 3K), the diameter of rotational
256 trajectory varied across different ΔQ states in the trained networks, suggesting that additional
257 biological constraints that were not considered for RNN training may have imposed a constant
258 diameter in the mouse brain. In addition to the analysis of ΔQ estimates from a RL model fit to
259 behaviors, we examined the coding persistency of the ground truth ΔQ which is available as the
260 activity of the action output neuron in each RNN agent. We confirmed that the ground truth ΔQ
261 was also persistently encoded in both optimal and mouse-like RNN agents (Figure S6).

262

263 **Persistency facilitates reliable and robust value retrieval by downstream neural networks**
264 The emergence of ΔQ coding persistency in RNN agents suggests that persistent coding is a
265 preferred solution in the task. What would be the advantage of persistent coding? One
266 possibility is that untangled persistency may allow a more reliable signal retrieval by the
267 downstream network to guide the action selection. We tested this possibility by training artificial
268 RNNs to retrieve the ΔQ signal from different temporal patterns of simulated population activity
269 (Figure 6A). For this purpose, RNNs are biologically relevant as they receive time-varying inputs
270 sequentially, as opposed to other decoder models (e.g. regression models).

271 We created artificial population activity encoding ΔQ in 4 different patterns: persistent,
272 and 3 types of non-persistent coding (Figure 6B). In persistent coding, 20% of cells encode ΔQ
273 as rate coding persistently. The slope of ΔQ tuning curve for each neuron was taken from its
274 distribution among RSC neurons (Figure S7). For the first 2 types of non-persistent coding, the
275 cellular identity of the persistent coding pattern was shuffled independently at each time bin to
276 alter the ΔQ persistency of each neuron without altering the population-level ΔQ signal in each
277 time bin. Non-persistent 1 allowed each neuron to encode ΔQ in multiple time points, while non-
278 persistent 2 was constrained that each neuron encodes ΔQ in only one of the 5 time points. In
279 the third non-persistent coding scheme, binary signals (active or inactive) at each time bin were
280 used to encode ΔQ by activating distinct sequences of neurons across time for different values
281 of ΔQ . We prepared 10 different sequences for 10 bins of ΔQ values.

282 Using these activity patterns as inputs, we trained RNNs to retrieve ΔQ . Various levels of
283 noise were added to the input activity to test a range of signal-to-noise ratio (SNR). The RNN
284 trained with the persistent ΔQ codes was able to retrieve ΔQ better than those trained with non-
285 persistent codes, especially when the input activity noise was high (Figure 6C-D). This indicates
286 that persistent coding facilitates reliable information retrieval by downstream circuits.
287 Furthermore, the RNNs that were trained to retrieve ΔQ from persistent coding were more
288 robust to changes in the synaptic weights, loss of synapses and cells (Figure 6E).

289 To investigate the impact of persistency in the brain activity, we next examined how ΔQ
290 could be retrieved from the neural activity with different levels of persistency recorded from the 6
291 cortical areas (Figure 6F). In addition to the original recorded activity ('Raw'), we artificially
292 increased or decreased ΔQ coding persistency by temporally sorting ('Sorted') or shuffling
293 ('Shuffled') the cell identity in each area. These persistency manipulations simply changed the
294 neuron ID of activity and thus did not alter the total amount of ΔQ signal in each time bin. Using
295 these sets of neural activity as inputs, we trained RNNs to retrieve ΔQ . There was a general
296 trend that an increase in persistency (*sorted activity*) improved retrieval accuracy, while a
297 decrease in persistency (*shuffled activity*) impaired retrieval accuracy (Figure 6G). However, the
298 effect size differed across different cortical areas. We found that the increase in retrieval
299 accuracy by sorting was larger when the original persistency in the population was lower, and
300 the decrease in retrieval accuracy by shuffling was larger when the original persistency was
301 higher (Figures 6H-I). These results further support the notion that coding persistency is a
302 critical determinant that enhances the accuracy of information retrieval by the downstream
303 network.

304 The results above indicate that persistent codes can be read out by the downstream
305 more effectively than non-persistent codes when the artificial neural network is allowed to train
306 its synaptic weights by minimizing the difference between its output and the target (ΔQ) as
307 supervised learning. However, in the real brain, such an explicit supervised target label to guide
308 the shaping of network connectivity is rarely available. Another approach to shape the
309 connectivity to retrieve particular information is unsupervised learning where errors are
310 computed using information readily available to the local network such as the input itself
311 (Lillicrap et al., 2020). Therefore, we next considered the possibility that coding persistency may
312 also affect signal retrieval processes that do not necessitate a supervised target label for each
313 information. It has been suggested that the brain may implement unsupervised learning in a
314 similar way to autoencoder networks in which the target is the input itself (Lillicrap et al., 2020).

315 Autoencoders extract the most dominant signals from the input activity and represent them in
316 the activity of a small number of neurons in the coding layer. The networks shape their
317 connectivity by reconstructing the input activity from the coding layer and minimizing the
318 reconstruction error between the input and the reconstructed activity. To examine what
319 information in the input population activity can be extracted in an unsupervised manner by
320 downstream recurrent networks, we used a recurrent denoising autoencoder (RDAE) (Maas et
321 al., 2012; Vincent et al., 2010) that sequentially processes input activity and extracts the latent
322 representations embedded in the input activity sequence, which are sufficient to reconstruct the
323 original population activity sequence with noise robustness (Figure 7B; Methods). When the
324 RDAE was trained on RSC population activity, ΔQ was extracted in the most dominant
325 dimensions of neural activity in the coding layer (Figure 7A). The ΔQ representation in the
326 coding layer was independent of upcoming choice directions, indicating that the dimensions
327 reflect value and not motor plans. Other task-related signals were not represented as the
328 dominant signals in the coding layer (Figure S8). Similar results were observed in the activity
329 dynamics of the mouse-like RNN agent but not in S1. Systematic comparisons among 6 cortical
330 areas revealed that extracted ΔQ in the coding layer was especially high from RSC, and the
331 amount of extracted ΔQ showed a high correlation with the ΔQ coding persistency in the input
332 population activity (Figures 7B-D). To directly test the effect of persistency, we artificially
333 manipulated the persistency of ΔQ coding in RSC without changing the total amount of ΔQ
334 signals in the population. We found that artificial increases in the persistency by sorting the cell
335 identity improved the ΔQ extraction, while artificial decreases in the persistency by shuffling the
336 cell identity worsened the ΔQ extraction (Figure 7E). These results indicate that high
337 persistency in the input activity can allow ΔQ retrieval by the downstream network even without
338 supervised learning.

339 Taken together, these analyses indicate that the persistency of value coding facilitates a
340 robust and accurate readout of value by downstream networks.

341

342 **Signal leakage can contribute to distributed value coding with varying levels of
343 persistency**

344 The results so far indicate computational advantages of persistent coding. However, in the
345 mouse brain, ΔQ coding was widely distributed across the 6 cortical areas with different levels
346 of persistency (Figures 1C-F). We asked whether anatomical connectivity among cortical areas
347 relates to the persistency levels of value coding. We analyzed the connectivity among imaged

348 areas using the dataset from the Allen Mouse Brain Connectivity Atlas (Oh et al., 2014).
349 Focusing on the projections from each of the 3 areas with high ΔQ persistency (RSC, PPC,
350 pM2), we quantified their axon projection density in each of the other 5 imaged areas (Figure
351 8A). We found that RSC, PPC, and pM2 predominantly project to each other, with smaller
352 amounts of direct projections to ALM, S1, and V1 (Figures 8B-C and S9). Thus, 3 areas with
353 persistent and strong ΔQ coding densely connect with each other, while they send less direct
354 projections to the other 3 areas with weaker and less persistent ΔQ coding. Based on this
355 observation, we hypothesized that the weak ΔQ persistency in ALM, S1 and V1 could result
356 from a signal leakage from the areas that maintain ΔQ as persistent activity. To test this
357 hypothesis, we built RNNs with multiple recurrent layers that receive RSC activity through non-
358 specific synaptic connectivity and examined how ΔQ coding changes along the downstream
359 hierarchy of layers (Figure 8D). We found that the fractions of neurons with ΔQ coding gradually
360 decreased as the signal leaked through layers of recurrent connectivity (Figures 8E-F).
361 Concurrently, ΔQ coding became increasingly less persistent (Figure 8G), and the temporal
362 tangling of ΔQ coding in neuronal manifolds gradually increased in the downstream (Figure 8H).
363 Furthermore, artificial manipulations of ΔQ coding persistency in the input RSC activity revealed
364 that persistency in ΔQ coding can affect the robust distribution of ΔQ coding with graded levels
365 of persistency across the downstream layers (Figures 8F-G). We obtained similar results using
366 PPC and pM2 as the input activity (Figure S10A-H), and the decreases in the ΔQ coding
367 neurons and the ΔQ coding persistency in the downstream layers were more dramatic when the
368 direct neural projections from layer to layer were sparse (Figure S10I-K). These results indicate
369 that, even without specific connectivity to selectively route particular information, persistently
370 encoded information can propagate thorough layers of non-specific connectivity to lead to a
371 wide distribution of the information encoded with lower levels of persistency in downstream
372 areas.

373

374 **DISCUSSION**

375 Brain-wide distribution of task-related information has emerged as a common principle in recent
376 years. In many cases, such as what we observed for ΔQ coding (Figure S1), task-related
377 signals are encoded by a heterogeneous population with some cells increasing but others
378 decreasing their activity. Such information coding may not be identified with classical large-scale
379 recording techniques such as fMRI, EEG and ECoG that quantify population average
380 responses. Even though information coding is wide-spread, the way by which information is

381 encoded differs across areas (Hattori et al., 2019). In the present study, the big data of >100k
382 mouse decisions and the activity from >100k neurons in 2 behavioral tasks allowed us to
383 investigate the potential origin of the distributed information coding and the computational
384 advantages of persistent coding using data-driven machine learning approaches. Coding
385 persistency was both learning- and context-dependent, and the persistent coding emerged
386 during task learning in both mouse brain and artificial network agents performing the same task.
387 Persistency facilitates an untangled maintenance of information as well as its reliable retrieval
388 by downstream circuits. The observation that persistency is context-dependent suggests that
389 certain cortical areas such as RSC can adjust coding persistency depending on behavioral
390 demands. For example, persistency may be especially preferred when the task context requires
391 extended maintenance of the information, or the maintained information is graded as in the case
392 of value, so that information can be stably maintained and robustly retrieved by downstream
393 areas. Furthermore, we showed that persistent coding in key areas such as RSC could also
394 contribute to the wide distribution of ΔQ coding across the mouse brain even through non-
395 specific signal leakage. The same principle may also apply to other task-related signals in
396 various task conditions, providing a possible explanation for the widespread phenomenon of
397 distributed coding across the brain. In other words, a wide distribution of information is expected
398 across the interconnected network of the brain, unless specific connectivity restricts the
399 propagation of particular information. We note that non-specific leakage is one of potential
400 mechanisms for signal distribution and it remains to be shown how much such a mechanism
401 contributes to the phenomenon. Furthermore, this mechanism is agnostic to whether the
402 propagated information has a function in the downstream areas — leaked information could
403 contribute to various computations performed in downstream areas.

404 We trained artificial RNNs to imitate the mouse behavioral strategy using behavioral
405 cloning and investigated the activity dynamics that emerged in the RNNs that were trained
406 without activity constraints. Previous studies trained task-performing artificial neural networks
407 either by using the correct action labels which are defined in each task structure (e.g. action A
408 must be taken after stimulus A) (Masse et al., 2019; Orhan and Ma, 2019) or by RL (Banino et
409 al., 2018; Song et al., 2017; Tsuda et al., 2020; Wang et al., 2018). Both approaches train the
410 networks to learn the optimal strategy in the respective task, independent of the actual
411 behavioral strategy that animals learn in the environment. In our value-based decision task,
412 animals learn to use behavioral history for decisions during training, but the RL strategy that
413 animals develop was suboptimal (Figure 4E-F). The origins of the sub-optimality likely include 1)
414 limited memory capacity, 2) low sample efficiency, 3) limited amount of training trials, and 4)

415 inductive bias inherent to each species. Deep RL, an artificial network that learns to solve a task
416 with RL, does not always have these constraints, and thus it learns a near-optimal strategy
417 unlike animals. These artificial networks may not reflect the mechanisms used by the brain. In
418 another common approach, simpler mathematical models (e.g. regression, classical RL models)
419 directly fit to animal behaviors are useful to understand the behavioral strategies. However, they
420 do not provide insights into potential neural activity dynamics that may mediate the behaviors.
421 To overcome these issues, we trained artificial RNNs, using mice as the teachers, to acquire the
422 sub-optimal RL strategy that mice develop during training. The big data of ~50k decisions
423 collected from expert mice allowed us to successfully train RNNs to imitate mouse behavioral
424 strategy. This data-driven approach to train RNNs to implement animal/human-like behaviors
425 would be a useful approach to obtain the neural network models and analyze what kind of
426 activity dynamics allows the animal strategy in a particular task. Similarly to our approach,
427 convolutional neural networks has been trained in visual object recognition tasks. The training
428 was done to perform the task optimally, as opposed to our approach using behavioral cloning.
429 Nevertheless these networks have been shown to develop some neural activity characteristics
430 that resemble the neural activity in the visual system of animals (Kriegeskorte, 2015; Yamins
431 and DiCarlo, 2016). These deep learning approaches will be a powerful approach to understand
432 what kind of neural activity may mediate given behaviors.

433 In this study, we developed dsPCA, a novel dimensionality reduction method which
434 combines the strengths of supervised and unsupervised algorithms. The supervised aspect
435 allows us to identify the best demixed linear coding dimensions for targeted task-related
436 variables, and the unsupervised aspect allows us to identify non-targeted correlated signals in
437 the remaining population activity. Therefore, dsPCA is a generally applicable method to
438 understand both the signals of interest and other non-targeted correlational structures in high-
439 dimensional data. Using dsPCA, we found that both mouse brain and artificial RNN agents
440 develop cylindrical dynamics, which consists of within-trial cyclic dynamics and its across-trial
441 transition along ΔQ axis. Similar within-trial dynamics have been well-studied in monkey motor
442 cortex during arm movement (Churchland et al., 2012; Russo et al., 2018, 2020). The studies
443 showed that the population activity state draws untangled rotational dynamics during
444 movements. They also showed that the activity state draws a simple cyclic trajectory in the
445 primary motor cortex, while the supplementary motor area draws a helical trajectory that unfolds
446 along a single direction by reflecting the 'context' of the movement (Russo et al., 2020). The
447 activity trajectory that we observed had cylindrical geometry, and the activity state repeatedly
448 transitioned along the ΔQ axis based on the RPE. These spatially confined geometries ensure

449 the untangled representation of ΔQ , which contributes to a robust ΔQ representation in the
450 brain. dsPCA and other RNN-based approaches in this study would facilitate the geometric
451 understanding of population dynamics in both biological and artificial networks.

452

453 **Acknowledgements**

454 We thank Marcus Benna and the Komiyama lab members for discussions and comments. This
455 work was supported by NIH (R01 NS091010, R01 EY025349, R01 DC014690, and P30
456 EY022589), NSF (1940181), and David & Lucile Packard Foundation to T.K., and the Uehara
457 Memorial Foundation Postdoctoral Fellowship, JSPS Postdoctoral Fellowship for Research
458 Abroad, and the Research Grant from the Kanae Foundation for the Promotion of Medical
459 Science to R.H.

460

461 **Author contributions**

462 R.H. and T.K. conceived the project. R.H. performed all the calcium imaging experiments,
463 analyses and simulations with suggestions from T.K. R.H. and T.K. wrote the paper.

464

465 **Declaration of Interests**

466 The authors declare no competing interests.

467

468 **Figure Legends**

469 **Figure 1. Persistency of action value coding across mouse cortex is area- and learning- 470 dependent.**

471 (A) Schematic of the value-based decision task and an example expert behavior.

472 (B) Neural activity was recorded from 6 cortical areas. The heatmap is the trial-averaged z-
473 scored deconvolved activity of an example RSC population. The activity of each neuron was
474 normalized to its peak. A half of the recorded trials were used to sort cells by the peak time, and
475 the mean activity of the other half are shown.

476 (C) Fractions of cells with significant ΔQ coding during ready period based on the mean activity
477 within each of the non-overlapping 200 ms bins (Regression, $P < 0.05$, 2-sided t-test). The

478 fractions with filled circles are significantly above the chance fraction of 5% ($P < 0.05$, one-sided
479 t-test). ΔQ values were shuffled across trials for the right panel.

480 (D) Activity of ΔQ coding neurons that were identified at different time windows (yellow
481 shadings) for example RSC and S1 populations. Trials were binned according to the ΔQ of each
482 trial, and the activity in each trial bin was averaged.

483 (E) t-values for ΔQ coding at each time bin of ready period for example populations of RSC and
484 S1 (Regression). Neurons were sorted based on the t-values at the last time bin.

485 (F) ΔQ coding persistency of each population as quantified by the persistency index (0: chance
486 persistency, 1: maximum-possible persistency, Methods, $***P < 0.001$, $****P < 0.0001$, one-way
487 ANOVA with Tukey's HSD).

488 (G) Fraction of trials when mice chose the side with higher reward assignment probability across
489 training sessions ($n = 9$ mice, mean \pm CI). The first 6 sessions were treated as early sessions.

490 (H) Activity of ΔQ coding RSC neurons that were identified from the activity within the specified
491 time bin (yellow shadings) in early and late sessions (same RSC population between the 2
492 sessions) indicating an increase in persistency during learning.

493 (I) ΔQ coding persistency of each population as quantified by the persistency index for early and
494 expert sessions ($**P < 0.01$, $***P < 0.001$, mixed effects model with population as the fixed
495 intercept).

496 All error bars are s.e.m.

497

498 **Figure 2. Persistency of history coding is task-dependent.**

499 (A) Schematic of the alternate choice task. The choice opposite to the choice in the previous
500 trial was rewarded regardless of reward outcome in the previous trial.

501 (B) Fraction of trials of correctly choosing the side with reward across training sessions ($n = 9$
502 mice, mean \pm 95% CI).

503 (C) Activity of RSC neurons that significantly encoded the action history from previous trial in
504 alternate choice task and value-based decision task. These neurons were identified using the
505 activity within the specified time bin (yellow shadings). The activity of the identified action history
506 coding neurons was separately averaged according to the choice on the previous trial.

507 (D) Persistency of action history coding in each population as quantified by the persistency
508 index for the alternate choice task (Alt) and the value-based decision task (Value) ($*P < 0.05$,
509 $**P < 0.01$, $***P < 0.001$, $****P < 0.0001$, mixed effects model with population as the fixed
510 intercept).

511

512 **Figure 3. dsPCA reveals cylindrical dynamics with untangled value representation in**
513 **RSC.**

514 (A) dsPCA decomposes the activity of a population of individual neurons that exhibit mixed
515 selectivity for multiple variables into demixed dimensions and the remaining subspace that is
516 free of the targeted signals.

517 (B) Matrix operations to identify the target-free axes. Full QR decomposition of a matrix with
518 target axes (\mathbf{T}) identifies a set of basis vectors that spans the target-free subspace (\mathbf{S}_{free}).

519 These target-free axes are realigned based on the principal component vectors ($\mathbf{W}_p^{\text{pca}}$, matrix
520 with top p PCA loadings) of the activity in the target-free subspace. The target-free axes in the
521 original n -dimensional space are the columns of $\mathbf{W}_p^{\text{dspca}} = \mathbf{S}_{\text{free}} \mathbf{W}_p^{\text{pca}}$.

522 (C) Fraction of activity variance along each target axis and the top 5 PC axes from the target-
523 free subspace. dsPCA was performed on noisy simulated data with target signals A, B, and C
524 (10 repeated simulations). The amount of variance is similar between the 5 target-free axes
525 because only Gaussian noise remained in the target-free subspace.

526 (D) Signals A, B, and C along each dimension identified with dsPCA for the simulated data.
527 Pearson correlations between the projected activity and each signal are shown.

528 (E) Decoding accuracy of target signals from original population activity, activity in the target
529 subspace (3 dimensions), and activity in the target-free subspace ($n-3$ dimensions). 50,000 and
530 10,000 trials for training and test sets.

531 (F) We applied dsPCA to decompose the original population activity into the demixed Q
532 subspace that consists of ΔQ , Q_{ch} , and ΣQ dimensions, and the Q-free subspace which is
533 orthogonal to the Q subspace.

534 (G) Fraction of activity variance along each Q-related axis and the top 5 PC axes from the Q-
535 free subspace for RSC populations. Unlike simulated data (C), the amount of variance between
536 axes of the Q-free subspace differ, indicating that non-targeted correlated signals exist in the Q-
537 free subspace.

538 (H) Q-related signals along each dsPCA dimension for RSC populations. Pearson correlations
539 between the projected activity and each signal are shown.

540 (I) Decoding accuracy of Q signals from the original RSC population activity, activity in the Q
541 subspace (3 dimensions), and activity in the target-free subspace.

542 (J-K) Example RSC, S1, and V1 population activity dynamics in neuronal manifolds where ΔQ
543 axis is paired with Q_{ch} and ΣQ axes (J), or axes that reflect major within-trial temporal activity

544 variance of Q-free subspace (K). dsPCA was applied on the activity between -2 and -1 sec from
545 choice, and the activity between ± 4 sec from choice was projected onto the identified axes.
546 Circles indicate the choice time. Projected activity was temporally downsampled to non-
547 overlapping 200 ms bins.
548 (L) Activity state transitions along ΔQ axis according to the updated action values in RSC,
549 whereas S1 and V1 activity draw complex trajectories that lead to tangling in the geometry.
550 Post-action selection trajectory was separately averaged according to the sign of ΔQ update.
551 (M) Activity state transitions in (L) shown along the ΔQ axis.
552 (N) Population activity in RSC forms cylindrical dynamics where within-trial cyclic dynamics can
553 transition along ΔQ axis across trials according to the RPE, while in the other areas ΔQ
554 representation is tangled.
555 All error bars are 95% CI.

556

557 **Figure 4. Untangled persistency emerges in the artificial RNNs trained to perform**
558 **‘mouse-like’ RL.**

559 (A) Optimal RNN agent was trained by updating its synaptic weights to minimize the
560 discrepancy in decisions (cross-entropy error) between the teacher (optimal choice generator)
561 and the student (RNN).
562 (B) Behaviors of the trained optimal RNN agent in an example session. The agent ran the task
563 by itself using its recurrent activity dynamics to implement RL. The left choice probability of the
564 RNN agent was taken from its output neuron activity. Left (Q_L) and right (Q_R) action values were
565 estimated by fitting a RL model to the behaviors.
566 (C) Mouse-like RNN agent was trained by updating its synaptic weights to minimize the
567 discrepancy in decisions (cross-entropy error) between the teacher (expert mice) and the
568 student (RNN).
569 (D) Behaviors of the trained mouse-like RNN agent in an example session.
570 (E) Frequency of rewarded trials (*left*) and choice predictability by a RL model optimized to
571 describe expert mouse behaviors (*right*, 5-fold cross-validation). $n = 82$ sessions for mice, 500
572 sessions (5 trained networks, each ran 100 sessions of 500 trials/session) each for the optimal
573 and mouse-like RNN agents.
574 (F) Decision dependence on history from past 10 trials, quantified by a regression model
575 (Methods). RewC: rewarded choice, UnrC: unrewarded choice, C: outcome-independent choice
576 history. $n = 82$ sessions for mice, 5 sessions (5 trained networks, each ran 10,000 trials) each

577 for the optimal and mouse-like RNN agents. The regression weights were normalized by the
578 model accuracy. Error bars are 95% CI.
579 (G) Activity of ΔQ coding neurons that were identified using the activity at the highlighted time
580 bin (yellow shading, -1 time step before choice) in the recurrent layer of a trained mouse-like
581 RNN agent (*left*), and the t-values of ΔQ after choice for the activity in each time bin (*right*).
582 Each trial had 10 time steps, and 0 corresponds to the choice time. t-values were sorted based
583 on the last time step (+9). The t-values in RNNs are higher than in mice due to smaller amount
584 of activity noise. Error bars are s.e.m.

585

586 **Figure 5. Cylindrical dynamics emerges in mouse-like RNN agents and mice during**
587 **training.**

588 (A) Population activity dynamics of the recurrent layer of mouse-like RNN agents in neuronal
589 manifolds where ΔQ axis is paired with axes that reflect major within-trial temporal activity
590 variance in Q-free subspace. Agents at each training stage ran the task for 10,000 trials. dsPCA
591 was applied on the activity averaged between -5 and -1 time steps from choice, and the
592 population activity between ± 5 time steps from choice was projected onto the identified axes. 4
593 independently trained mouse-like RNN agents are shown. Circles indicate the choice time.
594 (B) Population activity dynamics of example RSC, PPC, pM2, and ALM populations in early and
595 expert sessions. The same population of neurons was longitudinally compared for each area.
596 dsPCA was applied on the activity averaged between -2 and -1 sec from choice, and the
597 population activity between ± 4 sec is visualized. Circles indicate the choice time.

598

599 **Figure 6. Persistency in value coding facilitates reliable and robust value retrieval by**
600 **downstream neural networks.**

601 (A) RNN (40 recurrent units) was trained to retrieve ΔQ from the input population activity
602 sequence with either persistent or non-persistent ΔQ coding.
603 (B) Artificial population activity with either persistent or non-persistent ΔQ coding in the 200-cell
604 sequence. 3 types of non-persistent mode were considered (2 rate coding, 1 binary coding;
605 Methods). In the rate coding populations, the color indicates the Pearson correlation between
606 the activity and ΔQ (20 % of neurons at each bin encode ΔQ). Example populations were
607 visualized by either clustering ΔQ -coding neurons at each time bin (*top*) or sorting neurons
608 based on the correlation at the last time bin (*bottom*). In the binary coding population, ΔQ is

609 encoded by a unique activity sequence across time for each bin of ΔQ values (ten evenly
610 spaced bins between ± 1). 20% of neurons at each time bin participate in each sequence. In the
611 example, cells are sorted for either sequence 1 or 2. Time bins that are active in both
612 sequences are colored black.
613 (C) Mean ΔQ retrieval accuracy by the downstream RNNs from populations with different coding
614 modes and varying SNR (10 simulations for each).
615 (D) The ΔQ retrieval accuracy at the 5th time step with different SNR in the input activity. The
616 purple dashed line indicates the median SNR of ΔQ coding in imaged RSC populations.
617 (E) Robustness of trained RNNs. Simulations were performed using artificial population activity
618 with SNR of 1. Noise to synaptic weights was given by Gaussian noise with the standard
619 deviation relative to the standard deviation of the weight distribution of each connection type.
620 Error bars in (D) and (E) are 95% CI.
621 (F) Artificial manipulations of ΔQ coding persistency illustrated in an example PPC population
622 during ready period. Error bars are s.e.m.
623 (G) ΔQ retrieval accuracy before and after the persistency manipulations (subsampled 240 cells
624 were used, **P < 0.01, ****P < 0.0001, one-way ANOVA with Tukey's HSD).
625 (H) Gain in retrieval accuracy by sorting correlates with the original ΔQ coding persistency.
626 (I) Loss in retrieval accuracy by shuffling correlates with the original ΔQ coding persistency.

627

628 **Figure 7. Persistency in value coding also facilitates unsupervised value retrieval by
629 downstream neural networks.**

630 (A) Representation of input population activity in the coding layer of denoising recurrent
631 autoencoder networks (RDAE). Each network was trained to extract major signals from example
632 populations of RSC, S1, or a trained mouse-like RNN agent (5,000 trials). Population activity
633 sequence during ready period was used as the input. Each data point corresponds to a trial,
634 with the colors indicating the ΔQ of the trial. Trials were separated according to the choice
635 directions in the upcoming answer period in the bottom 2 rows. The dominant signals extracted
636 in the activity of coding neurons (10 neurons) were visualized in 2 dimensions by
637 multidimensional scaling.
638 (B) RDAEs extract major signals of the input population activity into the activity of N neurons in
639 the coding layer by unsupervised learning.
640 (C) Decoding accuracy of ΔQ from the activity of N neurons in the coding layer. A simple
641 feedforward neural network (N neurons in the coding layer are connected to a single output

642 neuron with *tanh* activation function) was used to decode from the coding layer. Input
643 populations were subsampled 240 cells.
644 (D) Decoding accuracy of ΔQ from the activity of neurons in the coding layer ($N = 1$ and 10)
645 positively correlates with the ΔQ coding persistency of the input population activity.
646 (E) Artificial manipulations of ΔQ coding persistency in the input RSC population bi-directionally
647 alter the amount of extracted ΔQ signal in the coding layer.
648 All error bars are s.e.m.

649

650 **Figure 8. Non-specific signal leakage can contribute to widely distributed value coding**
651 **with graded persistency.**

652 (A) Injection coordinates for anterograde tracing virus. RSC (red, $n = 60$ experiments), PPC
653 (blue, $n = 9$), and pM2 (yellow, $n = 33$). Experiments with left hemisphere injections were
654 mirrored horizontally. Experiments with both WT mice and Cre-transgenic mice were included
655 (See Figure S9 for WT only). White squares indicate the imaging FOVs used for our neural
656 activity analyses.
657 (B) Mean projection density of axons from each source area. Black dots indicate the injection
658 coordinates.
659 (C) Connectivity matrix with the mean projection density from each source area to the 6 target
660 areas that we used for our neural activity analyses ($500\mu\text{m} \times 500\mu\text{m}$ white squares).
661 (D) RSC population activity sequences were processed through 5 recurrent layers with non-
662 specific connectivity. Connection probability from layer to layer was set to 20% (Other
663 probabilities in Figure S10).
664 (E) Fractions of ΔQ coding neurons at each of the 200 ms time bins during ready period
665 (Regression, $P < 0.05$, 2-sided t-test). Error bars are s.e.m.
666 (F) Mean fractions of ΔQ coding neurons at each layer during ready period. Fractions of time
667 bins within the ready period were averaged for each population. Artificial manipulations of ΔQ
668 coding persistency in RSC does not affect the fractions of ΔQ coding neurons in RSC, but affect
669 the fractions in the downstream.
670 (G) ΔQ coding persistency at each layer. Persistency progressively decreases in the
671 downstream. Artificial manipulations of ΔQ coding persistency affect the persistency in the
672 downstream. Error bars in (F) and (G) are 95% CI.
673 (H) Temporal dynamics of population activity states visualized with dsPCA applied at each

674 layer. Cylindrical dynamics gradually collapses into highly tangled dynamics in downstream
675 layers.

676

677

678 **STAR Methods**

679 **Resource availability**

680 **Lead Contact**

681 Further information and requests for resources should be directed to and will be fulfilled by the
682 lead contact, Takaki Komiyama (tkomiyama@ucsd.edu).

683 **Materials Availability**

684 This study did not generate new unique reagents.

685 **Data and code availability**

686 • Data reported in this paper are available from the lead contact upon reasonable request.
687 • dsPCA code has been deposited at Zenodo and is publicly available. The DOI and the link to
688 the latest code in the GitHub repository are listed in the key resource table.
689 • Any additional information required to reanalyze the data reported in this paper is available
690 from the lead contact upon request.

691

692 **Experimental model and subject**

693 **Animals**

694 The experimental data in the value-based decision task were first reported in ref. (Hattori et al.,
695 2019). The data in the alternate choice task were newly collected for the current study. Both
696 male and female mice were included in both datasets because we did not observe obvious sex-
697 dependent differences in their neural activity patterns. Mice were originally obtained from the
698 Jackson Laboratory (CaMKIIa-tTA: B6;CBA-Tg(Camk2a-tTA)1Mmay/J [JAX 003010]; tetO-
699 GCaMP6s: B6;DBA-Tg(tetO-GCaMP6s)2Niell/J [JAX 024742]). All mice (6 weeks or older) were
700 implanted with glass windows above their dorsal cortex for *in vivo* two-photon calcium imaging.
701 All mice were water-restricted at ~1ml/day during training.

702

703 **Method details**

704 **Surgery**
705 Mice were continuously anesthetized with 1-2% isoflurane during surgery after subcutaneous
706 injection of dexamethasone (2mg/kg). After exposing the dorsal skull and removing the
707 connective tissue on the skull surface using a razor blade, we marked on the skull with black ink
708 at the coordinates of [AP from bregma, ML from bregma] = [+3.0 mm, 0 mm], [+2.0 mm, 0 mm],
709 [+1.0 mm, 0 mm], [0 mm, 0 mm], [-1.0 mm, 0 mm], [-2.0 mm, 0 mm], [-3.0 mm, 0 mm], [0 mm,
710 ±1.0 mm], [0 mm, ±2.0 mm], [0 mm, ±3.0 mm], [-2.0 mm, ±1.0 mm] , [-2.0 mm, ±2.0 mm], [-2.0
711 mm, ±3.0 mm]. We then applied saline on the skull and waited for a few minutes until the skull
712 became transparent enough to visualize vasculature patterns on the brain surface. We took a
713 photo of the vasculature patterns along with marked coordinates and used it to find target
714 cortical areas for two-photon microscopy. A large craniotomy was performed to expose 6
715 cortical areas, and a hexagonal glass window was implanted on the brain. The glass window
716 was secured on the edges of the remaining skull using 3M Vetbond (WPI), followed by
717 cyanoacrylate glue and dental acrylic cement (Lang Dental). After implanting the glass window,
718 a custom-built metal head-bar was secured on the skull above the cerebellum using
719 cyanoacrylate glue and dental cement. Mice were subcutaneously injected with Buprenorphine
720 (0.1 mg/kg) and Baytril (10 mg/kg) after surgery.

721
722 **Behavior task and training**
723 Mice were water-restricted at 1-2 ml/day after a minimum of 5 days of recovery after surgery.
724 We began animal training in pre-training tasks after at least a week of water restriction. We used
725 BControl (C Brody), a real-time system running on Linux communicating with MATLAB, to
726 control behavioral apparatus. We placed 2 lickports in front of head-fixed mice to monitor their
727 licking behaviors and give water rewards. Licking behaviors were monitored by IR beams
728 running in front of each water tube. We used an amber LED (5mm diameter) as the ready cue
729 and a speaker for auditory cues. Each trial begins with a ready period (2 or 2.5 sec with the
730 amber LED light), followed by an answer period with an auditory go cue (10 kHz tone). The 10
731 kHz tone was terminated when animals made a choice (the first lick to a lickport) or when the
732 answer period reached the maximum duration of 2 sec. Mice received a 50 ms feedback tone
733 (left: 5 kHz, right: 15 kHz) after a choice. ~2.5 µl water was provided to mice on each rewarded
734 trial from a lickport.

735 Before running in the alternate choice task or value-based decision task, mice were
736 trained in 2 pre-training tasks. In the 1st pre-training task, mice were rewarded for either choice
737 during the answer period. We gradually increased the mean ITI from 1 sec to 6 sec with ±1 sec

738 jitter. Through training in this task (2-3 days), mice learn that they can obtain water rewards from
739 the 2 lickports if they lick during the answer period. In the 2nd pre-training task, reward location
740 alternated every trial irrespective of their choice directions. Furthermore, licking during ready
741 period was punished by 500 ms white noise alarm sound and trial abort with an extra 2 sec ITI
742 in addition to the regular 5-7 sec ITI. Through training in this 2nd pre-training task (2-3 days),
743 mice learned to lick from both lickports and withhold licking during the ready period.

744 ***Alternate choice task***

745 In the alternate choice task, mice need to change their choice from a previous trial to get a
746 water reward. For example, if a mouse chose left on one trial, regardless of whether the mouse
747 received a reward or not, a water reward is available only from the right choice on the next trial.
748 The mouse will not get any rewards by repeating left choices for many trials because a reward
749 will not be assigned to the left until the mouse collects the assigned reward on the right side.
750 Mice need to rely on which side they chose in the previous trial to make the correct choice. ITI
751 was 5-7 sec, and the trials with licking during ready period were classified as alarm trials (500
752 ms white noise alarm sound and extra 2 sec ITI). Mice were trained for at least 2 weeks before
753 starting 2-photon calcium imaging.

754 ***Value-based decision task***

755 In the value-based decision task, a reward is probabilistically assigned to each choice. On each
756 trial, a reward may be assigned to each choice according to the reward assignment probabilities
757 that are different between two choices. Once a reward was assigned to a lickport, the reward
758 remained assigned until it was chosen. As a result, the probability that a reward is assigned to a
759 choice gradually increases if the choice has not been selected in the recent past trials. The
760 combinations of reward assignment probabilities were either [60 %, 10 %] or [52.5 %, 17.5 %] in
761 a trial, and reward assignment probabilities switched randomly every 60-80 trials in the order of
762 [Left, Right] = ..., [60 %, 10 %], [10 %, 60 %], [52.5 %, 17.5 %], [17.5 %, 52.5 %], [60 %, 10 %],
763 The probability switch was postponed if the fraction of choosing the lickport with higher
764 reward assignment probability was below 50 % in recent 60 trials until the fraction reached at
765 least 50 %. ITI was 5-7 sec, and the trials with licking during ready period were classified as
766 alarm trials (500 ms white noise alarm sound and extra 2 sec ITI). Trials in which mice licked
767 during ready period ('alarm trials', 5.15 %) and the trials in which mice failed to lick during the
768 answer period ('miss trials', 4.68 %) were not rewarded. We did not include alarm and miss
769 trials in neural activity analyses to ensure that the ready periods we analyzed were free of
770 licking behaviors and that mice were engaged in the task in the trials.

772 **Two-photon calcium imaging**

773 We used a two-photon microscope (B-SCOPE, Thorlabs) with a 16 \times objective (0.8 NA, Nikon)
774 and 925 nm excitation wavelength (Ti-Sapphire laser, Newport) for *in vivo* calcium imaging.
775 Images were acquired using ScanImage (Vidrio Technologies) running on MATLAB. All calcium
776 imaging was performed using camk2-tTA::tetO-GCaMP6s double transgenic mice that express
777 GCaMP6s in camk2-positive excitatory neurons. Each field-of-view (FOV) (512 \times 512 pixels
778 covering 524 \times 524 μ m) was scanned at \sim 29 Hz. Areas within the FOV that were not
779 consistently imaged across frames were discarded from analyses (Typically 10 pixels from each
780 edge of the FOV). We imaged and analyzed layer 2/3 neurons of 6 cortical areas in this study:
781 retrosplenial (RSC, 0.4 mm lateral and 2 mm posterior to bregma), posterior parietal (PPC, 1.7
782 mm lateral and 2 mm posterior to bregma), posterior premotor (pM2, 0.4 mm lateral and 0.5 mm
783 anterior to bregma), anterior lateral motor (ALM, 1.7 mm lateral and 2.25 mm anterior to
784 bregma), primary somatosensory (S1, 1.8 mm lateral and 0.75 mm posterior to bregma), and
785 primary visual (V1, 2.5 mm lateral and 3.25 mm posterior to bregma) cortex. Images from these
786 areas were collected from both hemispheres. We collected only 1 population from each
787 hemisphere for each cortical area of a single mouse. We imaged both hemispheres in two
788 different behavioral sessions if the FOVs on both hemispheres were clear at the time of
789 imaging.

790

791 **Image processing**

792 Images from 2-photon calcium imaging were processed using a custom-written pipeline (Hattori,
793 2021). The pipeline corrects motion artifacts using pyramid registration (Mitani and Komiyama,
794 2018), and slow image distortions were further corrected by affine transformations based on
795 enhanced correlation coefficients between frames (Evangelidis and Psarakis, 2008). We used
796 Suite2P (Pachitariu et al., 2016) to define regions of interests (ROIs) corresponding to individual
797 neurons and extract their GCaMP fluorescence. We selected only cellular ROIs using a user-
798 trained classifier in Suite2P and by manual inspections. At the step of signal extraction from
799 each cellular ROI, we excluded pixels that overlap with the other ROIs.

800

801 **Neural activity**

802 The neural activity data for the value-based decision task were first reported in ref. (Hattori et
803 al., 2019). We also additionally collected new neural activity data from mice running the
804 alternate choice task. The activity was continuously recorded with *in vivo* two-photon calcium
805 imaging at \sim 29 Hz from mice during the task performance. GCaMP fluorescence time series

806 were deconvolved to obtain signals that better reflect the kinetics of neural spiking activity using
807 a non-negative deconvolution algorithm (Friedrich et al., 2017; Pachitariu et al., 2018). The
808 deconvolved signal of each neuron was z-score normalized using the activity time series during
809 the entire imaging session before performing all the activity analyses in this study.

810 For the alternate choice task, we collected and analyzed the activity of 8,524 RSC
811 neurons (14 populations), 3,186 PPC neurons (7 populations), 7,915 pM2 neurons (14
812 populations) and 4,911 ALM neurons (10 populations) from 9 expert mice while they were
813 running the alternate choice task. For the value-based decision task, we analyzed the activity of
814 9,254 RSC neurons (15 populations), 6,210 PPC neurons (13 populations), 7,232 pM2 neurons
815 (13 populations) and 5,498 ALM neurons (10 populations) from early sessions ($\leq 6^{\text{th}}$ session),
816 and 9,992 RSC neurons (populations), 7,703 PPC neurons (populations), 9,759 pM2 neurons (
817 populations), 6,721 ALM neurons (populations), 7,576 S1 neurons (14 populations) and 2,767
818 V1 neurons (6 populations) from expert sessions of the data used in ref. (Hattori et al., 2019).
819

820 **Reinforcement learning model for mouse behaviors**

821 The reinforcement learning model that we used to estimate the action values in each trial was
822 taken from ref. (Hattori et al., 2019). This model was optimized specifically for mouse behaviors
823 and not necessarily ideal for describing the RL action policy of artificial neural network agents
824 (e.g. Optimal RNN agents). Action values of chosen (Q_{ch}) and unchosen (Q_{unch}) options in each
825 trial were updated as follows:

$$826 Q_{ch}(t+1) = \begin{cases} Q_{ch}(t) + \alpha_{rew} * (R(t) - Q_{ch}(t)) & \text{if rewarded } (R(t) = 1) \\ Q_{ch}(t) + \alpha_{unr} * (R(t) - Q_{ch}(t)) & \text{if unrewarded } (R(t) = 0) \end{cases} \quad [\text{eq. 1}]$$

$$827 Q_{unch}(t+1) = (1 - \delta) * Q_{unch}(t) \quad [\text{eq. 2}]$$

828 where α_{rew} and α_{unr} are the learning rates for rewarded and unrewarded trials respectively, δ is
829 the forgetting rate for the unchosen option, and $R(t)$ is reward outcome in trial t (1 for rewarded,
830 0 for unrewarded trials). The learning rates and the forgetting rate were constrained between 0
831 and 1. In alarm and miss trials, values of both options were discounted by δ . The probability of
832 choosing left (P_L) on trial t is estimated using left (Q_L) and right (Q_R) action values as follows:

$$833 P_L(t) = \frac{1}{1 + e^{-\beta_{\Delta Q}(\beta_0 + Q_L(t) - Q_R(t))}} \quad [\text{eq. 3}]$$

834 where β_0 is the value bias which is constant within each session, and $\beta_{\Delta Q}$ reflects the behavioral
835 sensitivity to ΔQ . The RL model was fit to the behavioral choice patterns with maximum
836 likelihood estimation.

837

838 **ΔQ-coding neurons**

839 ΔQ -coding neurons in the value-based decision task were identified with the following multiple
840 linear regression model.

841
$$a_i(t) = \beta_C C(t) + \beta_{\Delta Q} \Delta Q(t) + \beta_{Q_{ch}} Q_{ch}(t) + \beta_{\Sigma Q} \Sigma Q(t) + \beta_0 \quad [\text{eq. 4}]$$

842 where $a_i(t)$ is the mean activity of i^{th} neuron within each 200 ms time bin on trial t (except for
843 some analyses (Figures S1 and S2) where the mean activity within the first 2 sec of ready
844 period was used instead), $C(t)$ is the choice on trial t (1 if contralateral choice, -1 if ipsilateral
845 choice), $\Delta Q(t)$ is the value difference between contralateral and ipsilateral options on trial t ,
846 $Q_{ch}(t)$ is the value of the chosen option on trial t , and $\Sigma Q(t)$ is the sum of values of both
847 options on trial t . The regression weights were estimated by the ordinary least squares method.
848 ΔQ -coding neurons were identified with two-tailed t-test for the $\beta_{\Delta Q}$ regression weight (statistical
849 threshold of either $P < 0.05$ or $P < 0.01$ as indicated in the figure legend of each analysis). The
850 t-value for $\beta_{\Delta Q(t)}$ is $T_{\beta_{\Delta Q(t)}} = \frac{\beta_{\Delta Q}}{se(\beta_{\Delta Q})}$ where $se(\beta_{\Delta Q})$ is an estimate of the standard error of $\beta_{\Delta Q}$.

851

852 **Action history coding neurons**

853 Neurons that encode action history from an immediately preceding trial in the alternate choice
854 task and the value-based decision task were identified with the following multiple linear
855 regression model.

863
$$a_i(t) = \beta_{C_t} C(t) + \beta_{C_{(t-1)}} C(t-1) + \beta_0 \quad [\text{eq. 5}]$$

864 where $a_i(t)$ is the mean activity of i^{th} neuron within each 200 ms time bin on trial t , $C(t)$ is the
865 choice on trial t (1 if contralateral choice, -1 if ipsilateral choice), $C(t-1)$ is the choice on trial
866 ($t-1$) (1 if contralateral choice, -1 if ipsilateral choice, 0 otherwise). The regression weights
867 were estimated by the ordinary least squares method. Action history coding neurons were
868 identified with two-tailed t-test for the $\beta_{C_{(t-1)}}$ regression weight (statistical threshold of $P < 0.05$).

869 The t-value for $\beta_{C_{(t-1)}}$ is $T_{\beta_{C_{(t-1)}}} = \frac{\beta_{C_{(t-1)}}}{se(\beta_{C_{(t-1)}})}$ where $se(\beta_{C_{(t-1)}})$ is an estimate of the standard error
870 of $\beta_{C_{(t-1)}}$.

871

865 **Persistency index**

866 Persistency index to quantify the mean persistency of ΔQ coding or action history coding in a
 867 population of neurons was defined as follow;

$$868 \text{ Persistency index} = \frac{\frac{1}{m} \sum_{j=1}^m \sum_{i=1}^n \text{std}(T_{\text{shuffled sequence}}^{i,j}) - \sum_{i=1}^n \text{std}(T_{\text{raw sequence}}^i)}{\frac{1}{m} \sum_{j=1}^m \sum_{i=1}^n \text{std}(T_{\text{shuffled sequence}}^{i,j}) - \sum_{i=1}^n \text{std}(T_{\text{sorted sequence}}^i)} \quad [\text{eq. 6}]$$

869 where $T_{\text{raw sequence}}^i$ is the time series of t-values for $\beta_{\Delta Q}$ or $\beta_{C_{(t-1)}}$ that was obtained by fitting
 870 the [eq. 4] or [eq. 5] to the activity of each of the non-overlapping 200 ms time bins between 5
 871 sec before the ready cue and 2 sec after the ready cue. The across-time standard deviation of
 872 the $T_{\text{raw sequence}}^i$ was summed across all n neurons in the population (including neurons with
 873 non-significant t-values), and this summed standard deviation was normalized by min-max
 874 normalization such that the persistency index ranges between 0 (chance level persistency of a
 875 target population) and 1 (maximum persistency of a target population). The maximum
 876 persistency of a target population, $\sum_{i=1}^n \text{std}(T_{\text{sorted sequence}}^i)$, was obtained by independently
 877 sorting the cell identity at each time bin according to the $\beta_{\Delta Q}$ t-values of each cell in the time bin.

878 The chance level persistency of a target population, $\frac{1}{m} \sum_{j=1}^m \sum_{i=1}^n \text{std}(T_{\text{shuffled sequence}}^{i,j})$, was
 879 obtained by independently shuffling the cell identity at each time bin. To minimize the effect of
 880 randomness in the shuffling procedure, we iterated the shuffling m times ($m = 10$) and took the
 881 mean of the 10 iterations. This persistency index describes how persistent the target signal
 882 coding is above chance and how far the persistency is from the maximum persistency that the
 883 target population activity could achieve.

884

885 **Demixed subspace principal component analysis (dsPCA)**

886 Supervised dimensionality reduction algorithms can identify dimensions that encode targeted
 887 signals in high-dimensional data. However, they do not provide any information about signals
 888 that are not targeted by the users. As a result, these supervised analyses may miss important
 889 signals that exist in the original high-dimensional data. On the other hand, unsupervised
 890 dimensionality reduction algorithms can find dimensions for the major signals in the high-
 891 dimensional data, but they do not automatically reveal what kind of signals are reflected along
 892 each dimension. Furthermore, unsupervised methods may miss the signals of interest if the
 893 target signals are much weaker than the other dominant signals in the data.

894 We developed a novel dimensionality reduction algorithm that combines the strengths of
 895 both supervised and unsupervised methods. The demixed subspace principal component

896 analysis (dsPCA) identifies demixed coding axes for targeted variables in a supervised manner,
 897 and then identify axes that capture the remaining variance in the data using an unsupervised
 898 method. Although previously reported demixed principal component analysis (dPCA) has similar
 899 objectives (Kobak et al., 2016), dPCA can only identify targeted coding axes for discrete
 900 variables. In contrast, dsPCA can identify demixed axes for both discrete and continuous
 901 variables. Furthermore, although dPCA splits each targeted signal into multiple linear axes,
 902 dsPCA identifies a single linear coding dimension for each of the target signals, and all the
 903 linear information for the target signals are contained within the dimensions identified by these
 904 single coding axes.

905 The input to the algorithm is a 3rd-order tensor of population activity with dimensions of
 906 Trial (m) \times Time (t) \times Neuron (n).

907
$$X_{trial \times time \times neuron} = X_{m \times t \times n} \quad [eq. 7]$$

908 The tensor $X_{m \times t \times n}$ is first averaged over time axis elements within a specified time range, and
 909 we get a 2nd-order tensor of $X'_{m \times n}$.

910
$$X'_{m \times n} = \begin{pmatrix} x_{1,1} & x_{2,1} & \cdots & x_{n,1} \\ x_{1,2} & x_{2,2} & \cdots & x_{n,2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1,m} & x_{2,m} & \cdots & x_{n,m} \end{pmatrix} \quad [eq. 8]$$

911 To identify the demixed linear coding axes that encode ΔQ , Q_{ch} , or ΣQ in the population
 912 activity, we fit the following multiple linear regression model to the mean activity of individual
 913 neurons during the ready period;

914
$$\begin{pmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,m} \end{pmatrix} = \begin{pmatrix} 1 & \Delta Q_1 & Q_{ch_1} & \Sigma Q_1 \\ 1 & \Delta Q_2 & Q_{ch_2} & \Sigma Q_2 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \Delta Q_m & Q_{ch_m} & \Sigma Q_m \end{pmatrix} \begin{pmatrix} \beta_{i,0} \\ \beta_{i,\Delta Q} \\ \beta_{i,Q_{ch}} \\ \beta_{i,\Sigma Q} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{pmatrix} \quad [eq. 9]$$

915 where $\beta_{i,\Delta Q}$, $\beta_{i,Q_{ch}}$, and $\beta_{i,\Sigma Q}$ are the regression coefficients of the i^{th} neuron. For a population of
 916 n neurons, we obtain n regression coefficients for each type of Q-related signal. These
 917 regression coefficients are used to define the coding axes as follows;

918
$$\overrightarrow{\Delta q} = \begin{pmatrix} \Delta q_1 \\ \Delta q_2 \\ \vdots \\ \Delta q_n \end{pmatrix} = \frac{\overrightarrow{\beta_{\Delta Q}}}{\|\overrightarrow{\beta_{\Delta Q}}\|_2} = \frac{(\beta_{1,\Delta Q} \quad \beta_{2,\Delta Q} \quad \cdots \quad \beta_{n,\Delta Q})^T}{\sqrt{\sum_{i=1}^n |\beta_{i,\Delta Q}|^2}} \quad [eq. 10]$$

925
$$\overrightarrow{\mathbf{q}_{ch}} = \begin{pmatrix} q_{ch_1} \\ q_{ch_2} \\ \vdots \\ q_{ch_n} \end{pmatrix} = \frac{\overrightarrow{\boldsymbol{\beta}_{Q_{ch}}}}{\|\overrightarrow{\boldsymbol{\beta}_{Q_{ch}}}\|_2} = \frac{(\beta_{1,Q_{ch}} \ \beta_{2,Q_{ch}} \ \cdots \ \beta_{n,Q_{ch}})^T}{\sqrt{\sum_{i=1}^n |\beta_{i,Q_{ch}}|^2}} \quad [\text{eq. 11}]$$

926
$$\overrightarrow{\Sigma \mathbf{q}} = \begin{pmatrix} \Sigma q_1 \\ \Sigma q_2 \\ \vdots \\ \Sigma q_n \end{pmatrix} = \frac{\overrightarrow{\boldsymbol{\beta}_{\Sigma Q}}}{\|\overrightarrow{\boldsymbol{\beta}_{\Sigma Q}}\|_2} = \frac{(\beta_{1,\Sigma Q} \ \beta_{2,\Sigma Q} \ \cdots \ \beta_{n,\Sigma Q})^T}{\sqrt{\sum_{i=1}^n |\beta_{i,\Sigma Q}|^2}} \quad [\text{eq. 12}]$$

919 Note that these coding axes are ‘demixed’ coding axes where the activity variance for partially
 920 correlated variables are demixed into one of the axes for the partially correlated variables
 921 thanks to the linear demixing in the regression model ([eq. 9]). Although some previous studies
 922 further orthogonalized these demixed coding axes (Mante et al., 2013), we did not orthogonalize
 923 between the coding axes because further orthogonalization would remix these best demixed
 924 coding axes.

927 Next, our goal is to identify a neural subspace that does not encode any of the targeted
 928 Q-related signals. To identify the neural subspace that is free of the 3 targeted Q-related
 929 signals, we solve the following full QR decomposition of an $n \times 3$ matrix with the 3 coding axis
 930 vectors using Householder reflections;

931
$$\begin{pmatrix} \Delta q_1 & q_{ch_1} & \Sigma q_1 \\ \Delta q_2 & q_{ch_2} & \Sigma q_2 \\ \vdots & \vdots & \vdots \\ \Delta q_n & q_{ch_n} & \Sigma q_n \end{pmatrix} = (\mathbf{S}_Q, \mathbf{S}_{\text{free}}) \mathbf{R}$$

932
$$= (\overrightarrow{\mathbf{q}_1}, \overrightarrow{\mathbf{q}_2}, \overrightarrow{\mathbf{q}_3}, \overrightarrow{\mathbf{f}_1}, \overrightarrow{\mathbf{f}_2}, \cdots \overrightarrow{\mathbf{f}_{(n-3)}}) \mathbf{R}$$

933
$$= \begin{pmatrix} q_{1,1} & q_{2,1} & q_{3,1} & f_{1,1} & f_{2,1} & \cdots & f_{(n-3),1} \\ q_{1,2} & q_{2,2} & q_{3,2} & f_{1,2} & f_{2,2} & \cdots & f_{(n-3),2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ q_{1,n} & q_{2,n} & q_{3,n} & f_{1,n} & f_{2,n} & \cdots & f_{(n-3),n} \end{pmatrix} \mathbf{R} \quad [\text{eq. 13}]$$

934 where \mathbf{R} is an upper triangular matrix, \mathbf{S}_Q is a neural subspace that captures all Q-related
 935 signals, and \mathbf{S}_{free} is the Q-free subspace that is orthogonal to the \mathbf{S}_Q . \mathbf{S}_Q is formed by 3
 936 orthonormal basis vectors $(\overrightarrow{\mathbf{q}_1}, \overrightarrow{\mathbf{q}_2}, \overrightarrow{\mathbf{q}_3})$, and these basis vectors and the 3 coding axis vectors
 937 $(\overrightarrow{\Delta \mathbf{q}}, \overrightarrow{\mathbf{q}_{ch}}, \overrightarrow{\Sigma \mathbf{q}})$ span the identical neural subspace. On the other hand, \mathbf{S}_{free} is formed by $(n-3)$
 938 target-free orthonormal vectors $(\overrightarrow{\mathbf{f}_1}, \overrightarrow{\mathbf{f}_2}, \cdots \overrightarrow{\mathbf{f}_{(n-3)}})$ and capture all the remaining population
 939 activity variance that were not captured by the subspace \mathbf{S}_Q . The representation of the
 940 population activity $\mathbf{X}'_{m \times n}$ in \mathbf{S}_{free} is given by

941
$$\text{proj}_{\mathbf{S}_{\text{free}}} \mathbf{X}' = \mathbf{X}' \mathbf{S}_{\text{free}} \quad [\text{eq. 14}]$$

942 Lastly, we further realign the dimensions of the Q-free subspace S_{free} such that
943 minimum numbers of dimensions are necessary to explain the remained activity variance as
944 much as possible. This realignment is done using the principal component vectors from PCA on
945 $\text{proj}_{S_{\text{free}}} X'$. The top p principal component vectors ($p \leq n - 3$) can be used as the major Q-free
946 subspace dimensions for dimensionality reduction purpose as follows;

947
$$F_p' = X' (S_{\text{free}} W_p^{\text{pca}}) = X' W_p^{\text{dspca}} \quad [\text{eq. 15}]$$

948 where the m -by- p matrix F_p' is the top p principal components of the activity within the Q-free
949 subspace, the $(n-3)$ -by- p matrix W_p^{pca} is the loadings matrix of the PCA, and the n -by- p matrix
950 W_p^{dspca} is the loadings matrix of the dsPCA. The columns of W_p^{dspca} are the Q-free axis
951 vectors in the raw n -dimensional population activity space. More generally, the neural subspace
952 that is free of k targeted variables can be obtained by the same [eq. 15] with $p \leq n - k$.

953 Through these steps ([eq. 7] ~ [eq.15]), dsPCA identified the 3 linearly demixed coding
954 axes for the targeted Q-related signals ($\vec{\Delta q}$, $\vec{q_{ch}}$, $\vec{\Sigma q}$), and ($n - 3$) target-free axes (column
955 vectors of W_{n-3}^{dspca}). We confirmed that none of the targeted signals could be linearly
956 decodable from the population activity within the obtained target-free subspace (Figures 3E, 3I
957 and S3C).

958 In this manuscript, we decomposed neural population activity into demixed Q subspace
959 and Q-free subspace using dsPCA. The Q subspace consists of demixed linear coding axes for
960 ΔQ , Q_{ch} and ΣQ , and all activity variance that linearly relates to these Q-related signals are
961 included in this subspace. On the other hand, all the other activity variance that did not remain
962 in the Q subspace is included in the Q-free subspace. The activity state of the neural population
963 changes across trials within the Q subspace depending on how each of the Q-related signals is
964 updated by choice and its outcome. We also identified the axes that capture the major within-
965 trial temporal activity variance in the Q-free subspace by performing PCA on the 2nd-order
966 tensors that are obtained by averaging $\text{proj}_{S_{\text{free}}} X$ over trial axis elements.

967

968 **Quantification of Q-related signals in subspaces from dsPCA**

969 dsPCA decomposed population activity into Q subspace and Q-free subspaces. We examined
970 the amount of Q-related signals in each subspace. The strength of Q-related signals in a full
971 population activity with n neurons was quantified using linear decoders given by

972
$$\Delta Q(t) = \sum_{i=1}^n \beta_i^{\Delta Q} a_i(t) + \beta_0^{\Delta Q} \quad [\text{eq. 16}]$$

973
$$Q_{ch}(t) = \sum_{i=1}^n \beta_i^{Q_{ch}} a_i(t) + \beta_0^{Q_{ch}} \quad [\text{eq. 17}]$$

974
$$\Sigma Q(t) = \sum_{i=1}^n \beta_i^{\Sigma Q} a_i(t) + \beta_0^{\Sigma Q} \quad [\text{eq. 18}]$$

975 where $a_i(t)$ is the activity of the i^{th} neuron on trial t , β_i^x is the regression weight for $a_i(t)$, and β_0^x
 976 is the constant term. The decoder was trained with an L2 penalty by selecting the regularization
 977 parameter by 5-fold cross-validation. The decoding accuracy was obtained with 5-fold cross-
 978 validation by separating trials into training and test sets. Similarly, the strength of Q-related
 979 signals in the 3-dimensional Q subspace and the $(n - 3)$ -dimensional Q-free subspaces were
 980 quantified using linear decoders on the projected population activity in each subspace as
 981 follows;

982
$$\Delta Q(t) = \sum_{i=1}^x \beta_i^{\Delta Q} s_i(t) + \beta_0^{\Delta Q} \quad [\text{eq. 19}]$$

983
$$Q_{ch}(t) = \sum_{i=1}^x \beta_i^{Q_{ch}} s_i(t) + \beta_0^{Q_{ch}} \quad [\text{eq. 20}]$$

984
$$\Sigma Q(t) = \sum_{i=1}^x \beta_i^{\Sigma Q} s_i(t) + \beta_0^{\Sigma Q} \quad [\text{eq. 21}]$$

985 where $s_i(t)$ is the population activity along the i^{th} dimension of the subspace on trial t , β_i^x is the
 986 regression weight for $s_i(t)$, and β_0^x is the constant term. $x = 3$ for Q subspace while $x = n - 3$
 987 for Q-free subspace. These analyses revealed that all Q-related signals were captured by the
 988 Q-subspace, while Q-related signals were completely absent in the Q-free subspace (Figures
 989 3E, 3I and S3C).

990

991 **RNN agents with optimal or mouse-like RL strategy**

992 The RNN agents trained to perform RL in this study consisted of 2 neurons in the input layer,
 993 100 neurons in the recurrent layer, and 1 neuron in the output layer. The agents were trained to
 994 perform RL in the same behavior task environment with 10 time steps per trial. The 2 input

995 neurons receive choice and reward outcome information only at the time step immediately after
996 choice, and the history of the choice outcome information was maintained through the recurrent
997 connectivity in the downstream recurrent layer. The sequence of activity fed into the input
998 neurons was given as vectors with either choice or reward history labels in their elements. The
999 elements that correspond to the time steps immediately after choice took 1 for left choice and -1
1000 for right choice in the choice history vector, and the elements took 1 for reward outcome and -1
1001 for no-reward outcome in the reward history vector. These elements took 0 in miss trials. The
1002 other elements of the vectors were all zeros. We sequentially fed 100 time steps of sequences
1003 into these input neurons, and the network training was done with unroll length of 100 time steps
1004 for backpropagation through time. The choice input neuron and reward input neuron connect
1005 with neurons in the recurrent layer. The neurons in the recurrent layer are connected with each
1006 other through recurrent connections, which allows each recurrent neuron to receive outputs of
1007 the previous time steps. The output of the recurrent layer is given by

1008
$$\mathbf{y}_{(t)} = \tanh(\mathbf{W}_x \mathbf{x}_{(t)} + \mathbf{W}_y \mathbf{y}_{(t-1)} + \mathbf{b}) \quad [\text{eq. 22}]$$

1009 where $\tanh(\cdot)$ is a hyperbolic tangent activation function of the form $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$, $\mathbf{x}_{(t)}$ is a 2
1010 $\times 1$ vector containing the choice and reward information from a previous time step, $\mathbf{y}_{(t-1)}$ is a
1011 100×1 vector containing the layer's outputs at time step t , \mathbf{W}_x is a 100×2 matrix containing
1012 the connection weights for the inputs of the current time step, \mathbf{W}_y is a 100×100 matrix
1013 containing the connection weights for the outputs of the previous time step, and \mathbf{b} is a 100×1
1014 vector containing each neuron's bias term. The recurrent neurons send their outputs to the
1015 output neuron. The output neuron calculates the probability of selecting left action in the trial
1016 with a sigmoid activation function of the form $\sigma(z) = \frac{1}{1+e^{-z}}$. The agent then selects an action for
1017 the trial probabilistically by following the choice probability from the output neuron. This 3-layer
1018 RNN agent was trained to perform either an optimal RL strategy or the RL strategy that mice
1019 develop after training using its recurrent activity dynamics.

1020 To train the RNNs to perform optimal RL in the task environment, we directly utilized the
1021 reward assignment rule of the task. In the value-based decision task, a reward is assigned to
1022 each choice according to the reward assignment probabilities of each choice on each trial. Once
1023 a reward was assigned to a lickport, the reward was maintained on the choice until it was
1024 chosen by the animal. As a result, the probability that a reward is assigned to a choice gradually
1025 increases if the choice has not been selected in the recent trials. The actual cumulative reward
1026 probabilities of left and right choices are given by

1038
$$P_L(t) = 1 - \prod_{x=t-N_R(t)}^t \{1 - A_L(x)\} \quad [\text{eq. 23}]$$

1039
$$P_R(t) = 1 - \prod_{x=t-N_L(t)}^t \{1 - A_R(x)\} \quad [\text{eq. 24}]$$

1027 where $A_c(x)$ is the reward assignment probability of choice c on trial x , $N_c(t)$ is the number of
 1028 successive c choices before trial t (e.g. $N_R(t) = 3$ when the choice on (t-4) was left and the
 1029 choices on (t-3), (t-2), (t-1) were right). Therefore, an optimal choice generator would select a
 1030 choice with higher cumulative reward probability on each trial as follows;

1040
$$\text{Optimal choice} = \operatorname{argmax}_c \{P_c(t)\} \quad [\text{eq. 25}]$$

1031 We used this optimal choice generator as the teacher to train RNNs to learn a near-optimal RL
 1032 strategy. Unlike the optimal choice generator that knows the exact reward assignment
 1033 probabilities ($A_c(x)$) and the reward assignment rule, the RNNs are agnostic to these hidden
 1034 variables. Therefore, our goal is to train the RNNs to use only the past choice and reward
 1035 history to make choices that are similar to the choices made by the optimal choice generator. To
 1036 train the RNNs to imitate the behaviors of the optimal choice generator, we calculated binary
 1037 cross-entropy as the loss function to be minimized. The cross-entropy is given by

1041
$$H_p = -\frac{1}{M} \sum_{i=1}^M (a_i^{\text{optimal}} \log(p_i^{\text{RNN}}) + (1 - a_i^{\text{optimal}}) \log(1 - p_i^{\text{RNN}})) \quad [\text{eq. 26}]$$

1042 where M is the total number of training trials, a_i^{optimal} is 1 or 0 when the optimal choice generator
 1043 selected left or right action on the i^{th} trial respectively, and p_i^{RNN} is the left choice probability of the
 1044 RNN agent from its output neuron.

1045 To train RNNs to perform mouse-like RL that is suboptimal in the task environment, we
 1046 used 50,472 decision making trials of expert mice in the task environment. We fed the choice
 1047 and reward history that expert mice experienced into the RNNs, and trained the RNNs to imitate
 1048 the choice patterns of expert mice. To do this, we calculated the binary cross-entropy as the
 1049 loss function to be minimized. The cross-entropy is given by

1050
$$H_p = -\frac{1}{M} \sum_{i=1}^M (a_i^{\text{mouse}} \log(p_i^{\text{RNN}}) + (1 - a_i^{\text{mouse}}) \log(1 - p_i^{\text{RNN}})) \quad [\text{eq. 27}]$$

1051 where M is the total number of training trials, a_i^{mouse} is 1 or 0 when the expert mouse selected
 1052 the left or right action in the i^{th} trial respectively, and p_i^{RNN} is the left choice probability of the RNN
 1053 agent from its output neuron.

1054 For the training of both the optimal RNN agents and mouse-like RNN agents, the cross-
 1055 entropy loss was calculated at variable time steps for each trial to reflect the temporal variability
 1056 of the timing of decision making in this task (variable ITI, variable ready-period, variable reaction
 1057 time), and all the synaptic weights of the RNN agent were trained with backpropagation through
 1058 time. The training was optimized using mini-batch gradient descent with Nesterov momentum
 1059 optimization (learning rate of 0.001 and momentum of 0.9, batch size of 128), and the training
 1060 was terminated when the loss for a validation set (1/5 of trials) stopped decreasing for the
 1061 consecutive 50 epochs as a form of regularization (Early stopping). The trained RNN agents ran
 1062 the task in a simulated environment with the length of 500 trials/session, and the RL behavioral
 1063 strategy in the simulated environment was quantified by a RL model optimized to describe expert
 1064 mouse behaviors [eq. 1 – 3] and a logistic regression model [eq. 28].

1065

1066 **Quantification of history-dependent behavioral strategy**

1067 The quantification of behavioral strategy for mice and RNN agents was performed with either a
 1068 RL model [eq. 1 – 3] or a logistic regression model [eq. 36]. The logistic regression model
 1069 predicts an action in each trial based on 3 types of history from the past 10 trials. The model is
 1070 given by

$$1071 \quad \text{logit}(P_L(t)) = \sum_{i=1}^{10} \beta_{RewC(t-i)} * \text{RewC}(t-i) + \sum_{i=1}^{10} \beta_{UnrC(t-i)} * \text{UnrC}(t-i) \\ 1072 \quad + \sum_{i=1}^{10} \beta_{C(t-i)} * C(t-i) + \beta_0 \quad [\text{eq. 28}]$$

1073 where $P_L(t)$ is the probability of choosing left on trial t , $\text{RewC}(t - i)$ is the rewarded choice
 1074 history on trial $t - i$ (1 if rewarded left choice, -1 if rewarded right choice, 0 otherwise),
 1075 $\text{UnrC}(t - i)$ is the unrewarded choice history on trial $t - i$ (1 if unrewarded left choice, -1 if
 1076 unrewarded right choice, 0 otherwise), $C(t - i)$ is the outcome-independent choice history on
 1077 trial $t - i$ (1 if left choice, -1 if right choice, 0 otherwise). $\beta_{RewC(t-i)}$, $\beta_{UnrC(t-i)}$, and $\beta_{C(t-i)}$ are
 1078 the raw regression weights of each history predictor, and β_0 is the history-independent constant
 1079 bias term. The sizes of these raw weights reflect the relative contribution of each history variable

1080 to decision making in a behavior session. However, the weight size does not reflect the absolute
1081 strength of the contribution to decision making because the strength of each history effect on
1082 decision making is determined by not only the regression weight but also the choice prediction
1083 accuracy of the regression model. Therefore, we normalized the regression weights by the
1084 choice predictability of the regression model as follows;

1085
$$\text{Normalized } \beta_x = \left(\frac{N_{choice}^{correct}}{N_{choice}^{all}} - 0.5 \right) * \frac{\beta_x}{\sum_{i=1}^{10} (|\beta_{RewC(t-i)}| + |\beta_{UnrC(t-i)}| + |\beta_{C(t-i)}|) + |\beta_0|} \quad [\text{eq. 29}]$$

1086 where N_{choice}^{all} is the number of choice trials in the session, and $N_{choice}^{correct}$ is the number of choice
1087 trials that were correctly predicted by the [eq. 36]. Each regression weight is divided by the sum
1088 of absolute values of all the regression weights before being multiplied by the choice prediction
1089 accuracy. This normalization turns raw regression weights to reflect the fraction of choice
1090 predictability by each of the history variable. These normalized weights are comparable across
1091 different behavior sessions or mice because they reflect the absolute strength of each history
1092 event on decision making. We used these normalized weights to compare the history
1093 dependence of expert mice and trained RNN agents for their decision making.

1094

1095 **Artificial population activity sequence**

1096 Artificial population activity sequences with either persistent or non-persistent rate coding of ΔQ
1097 were created based on the distributions of the tuning curves of ΔQ coding among RSC neurons.
1098 Each population consisted of 200 neurons with 5 time bins, and we assigned 20% of neurons at
1099 each time bin to encode ΔQ . The tuning curve slope of ΔQ coding of each RSC neuron ($\beta_{\Delta Q}$)
1100 was measured by fitting [eq. 4] to the activity during ready period. We defined across-trial
1101 standard deviation of $\beta_{\Delta Q}\Delta Q(t)$ from [eq. 4] as the signal standard deviation of ΔQ coding. To
1102 derive the noise standard deviation, we first subtracted $\beta_{\Delta Q}\Delta Q(t)$ from the ready period activity
1103 sequence of each trial. The residual ready period activity sequences were then concatenated
1104 across trials. The standard deviation of the concatenated activity sequence was defined as the
1105 noise standard deviation. The SNR of ΔQ coding was defined as the ratio of the signal standard
1106 deviation to the noise standard deviation. The tuning curve slope for each activity time bin was
1107 randomly sampled without replacement from the distributions of ΔQ -coding neurons. The ΔQ
1108 signal was linearly encoded at each time bin according to the sampled tuning curve slope, and
1109 additional Gaussian noise was added to the neural activity. The other non- ΔQ coding activity
1110 time bins simply exhibited Gaussian noise. We created populations with 3 different types of rate

1111 coding modes (*Persistent*, *Non-persistent 1*, *Non-persistent 2*). In the populations with
1112 *Persistent* mode, the identical 20% of neurons encoded ΔQ at all 5 time bins. In the populations
1113 with *Non-persistent 1* mode, we randomly selected 20% of neurons at each time bin as the ΔQ -
1114 coding neurons and allowed each neuron to encode ΔQ with different tuning curve slopes at
1115 different time bins. *Non-persistent 2* mode is similar to *Non-persistent 1*, except that each
1116 neuron in the population encoded ΔQ at only one of the time bins.

1117 In addition to the 3 rate coding schemes, we also considered a coding mode that
1118 encodes ΔQ as specific sequential activity patterns across cells in a population. In this 3rd non-
1119 persistent coding mode (*Non-persistent 3*), neural activity at each time bin can take only binary
1120 states (0: inactive, 1: active). Therefore, this population encodes ΔQ using only the identity of
1121 active cells. We encoded 10 different sequences in a population such that each sequence
1122 uniquely corresponds to one of the 10 binned ΔQ (-1 to 1 with binning of 0.2 width). For each
1123 sequence, we randomly assigned 20% of neurons at each time bin as active neurons with a
1124 constraint that each neuron can be active only at a single time step in a sequence. After
1125 encoding the 10 different sequences in a population, we added Gaussian noise to the activity of
1126 each neuron. We defined the SNR of this coding scheme as the ratio of the across-time
1127 standard deviation of the activity of a neuron to the standard deviation of its added Gaussian
1128 noise.

1129 **ΔQ retrieval by RNN**

1130 RNNs were trained to retrieve ΔQ information from the input population activity sequence. The
1131 RNN had 40 recurrent neurons with *tanh* activation functions and an output neuron with linear
1132 activation function. The network weights were updated by backpropagation through time with
1133 RMSprop to minimize mean-squared-error (MSE) between the network outputs and ΔQ values
1134 of the trials in a training set. The network training was terminated when the MSE of a validation
1135 set stopped decreasing for the consecutive 20 epochs as a form of regularization (Early
1136 stopping). For each training iteration, we used 20% of available trials as a test set to calculate
1137 the ΔQ retrieval accuracy by the trained network, and the remaining 80% of the trials were
1138 further split into validation set (10%) and training set (70%). We repeated the network training 5
1139 times by using different sets of trials as the test set such that we can obtain ΔQ predictions by
1140 the trained networks for all available trials in a cross-validated way. The ΔQ retrieval accuracy
1141 was calculated by comparing the ΔQ predictions to the true ΔQ from the RL model. For the ΔQ
1142 retrieval from cortical activity, we used only 240 cells as the inputs to match the number of cells
1143 across different cortical areas. For each neural population, we subsampled 240 cells in each

1144 iteration allowing repetitions with the smallest number of iterations to include every cell at least
1145 once for decoding, and the ΔQ retrieval accuracy from the iterations were averaged.

1146

1147 **Denoising recurrent autoencoder**

1148 Autoencoder is an artificial neural network that learns to extract efficient coding of its input
1149 without supervision. It consists of an encoder network and a decoder network, and they are
1150 sequentially connected through a coding layer with small number of neurons. In a trained
1151 autoencoder, the encoder extracts essential signals in the input into the coding layer, while the
1152 decoder tries to reconstruct the original input from activity in the coding layer. When the number
1153 of neurons in the coding layer is smaller than the dimensions of the input, only signals that are
1154 dominant in the input remains in the coding layer of a trained autoencoder network. Among
1155 various types of autoencoders, we used denoising recurrent autoencoders (Maas et al., 2012;
1156 Vincent et al., 2010) to extract dominant signals embedded in each population activity
1157 sequence. Although autoencoders with only feedforward connections or convolutional neural
1158 networks can also extract latent signals in a population activity sequence, we used recurrent
1159 neural networks that sequentially process the input activity because the neural networks in a
1160 brain also process input activity sequentially. Our goal is to understand whether such
1161 biologically relevant recurrent networks can extract signals from input activity without explicit
1162 teaching labels (i.e. unsupervised learning). The latent signals extracted by a recurrent
1163 autoencoder represent the latent signals from the perspective of a recurrent network that
1164 processes input activity sequentially through its recurrent connectiviy.

1165 The autoencoders that we used to visualize extracted dynamics from example
1166 populations (Figure 7A) consisted of 3 hidden layers with recurrent connectivity (1st: 50 neurons,
1167 2nd: 10 neurons, 3rd: 50 neurons), and the activity of all neurons in a population was used as the
1168 input to the autoencoder. On the other hand, the autoencoders that we used for quantitative
1169 across-area comparisons (Figure 7B-E) consisted of 3 hidden layers with recurrent connectivity
1170 (1st: 20 neurons, 2nd: N neurons, 3rd: 20 neurons) and processed input activity of subsampled
1171 240 cells. Note that 3 layers are the minimum number of layers that are required for an
1172 autoencoder network. All recurrent neurons in the hidden layers had *tanh* activation functions.
1173 All neurons except for the neurons in the middle hidden layer (coding layer) sent activity
1174 sequentially to the neurons in the next layer. However, the neurons in the coding layer sent only
1175 the activity at the last time step to the next hidden layer. The last-time-step activity is the result
1176 of the temporal integration of the original population activity sequence through recurrent

1177 connectivity, and the activity reflects the latent representations in the original population activity
1178 sequence. The hidden layers after the coding layer reconstructed the original population activity
1179 sequence from the latent representations in the coding layer. The network weights were
1180 updated by backpropagation through time with RMSprop to minimize mean-squared-error
1181 (MSE) between the original population activity sequence and the reconstructed population
1182 activity sequence. To ensure stable training of network weights, we clipped the gradients of
1183 network weights if their L2 norms were greater than 1 (Gradient clipping (Pascanu et al., 2012)).
1184 To add noise robustness to the autoencoders, we applied dropout (Hinton et al., 2012;
1185 Srivastava et al., 2014) to the connections between the input neurons and the neurons in the 1st
1186 hidden layer such that 50% of randomly selected connections are ablated at each training step.
1187 The network training was terminated when the MSE of a validation set (20% of trials for Figure
1188 7A, 10% of trials for Figure 7B-E) stopped decreasing for the consecutive 20 epochs as another
1189 form of regularization (Early stopping). The activity of the 10 coding neurons for Figure 7A were
1190 further reduced to 2 dimensions with multidimensional scaling to visualize the dominant
1191 population activity states. To quantify the strength of ΔQ signal in the activity of N coding
1192 neurons for Figures 5B-E, we performed decoding of ΔQ from the activity of N coding neurons
1193 using a simple feedforward neural network where all the N coding neurons are connected to an
1194 output neuron with *tanh* activation function. For each training iteration, we used 20% of available
1195 trials as a test set to calculate the ΔQ decoding accuracy by the trained network, and the
1196 remaining 80% of the trials were further split into validation set (10%) and training set (70%).
1197 We repeated the network training 5 times by using different sets of trials as the test set such
1198 that we can obtain ΔQ predictions by the trained networks for all available trials in a cross-
1199 validated way. For these ΔQ decoding analyses, we also matched the number of cells included
1200 in the inputs to the autoencoders across different decoding by subsampling 240 cells from the
1201 original population. For each neural population, we subsampled 240 cells in each iteration
1202 allowing repetitions with the smallest number of iterations to include every cell at least once for
1203 decoding, and the ΔQ decoding accuracy from the iterations were averaged.

1204

1205 **Deep RNN with non-specific connectivity**

1206 Neural networks with 5 recurrent layers were used to simulate how the input population activity
1207 transforms in the downstream recurrent layers when the synaptic weights are non-specific
1208 throughout the networks. Each recurrent layer had 1,000 neurons with *tanh* activation functions,
1209 and the 5 recurrent layers were sequentially connected through feedforward connections. All

1210 neurons of a recorded cortical population were directly connected to the 1st recurrent layer.
1211 Each neuron in a recurrent layer was connected with all the other neurons in the same layer, but
1212 we made the connections between successive layers sparse by setting the connection
1213 probability of a neuron to the neurons in the next layer to 1%, 5%, 10%, 20%, or 50%. The non-
1214 specific synaptic weights were randomly drawn from a uniform distribution on [-1, 1].

1215

1216 **Anatomical connectivity analyses**

1217 We analyzed neural projections from the areas with high ΔQ coding persistency (RSC, PPC,
1218 pM2) using the neural tracing data available in the Allen Mouse Brain Connectivity Atlas (Oh et
1219 al., 2014). These projection data were originally acquired by injecting adeno-associated virus
1220 (AAV) encoding EGFP into various target brain areas and scanning EGFP-labelled axons
1221 throughout the brain with high-throughput serial 2-photon tomography. We used their software
1222 development kit (SDK), *allensdk*, to access and process their data in Python.

1223 ***Dorsal view of the Allen Reference Atlas***

1224 Allen Reference Atlas is a high-resolution anatomical 3D reference atlas for the adult mouse
1225 brain. Different brain structures are colored differently in this atlas. All projection data in the
1226 Connectivity Atlas are registered to this reference atlas. We created a dorsal view of the Allen
1227 Reference Atlas to indicate the virus injection coordinates and cortical projection density in the
1228 dorsal cortex. First, we downloaded the 3D RGB-colored atlas at the resolution of 25 $\mu\text{m}/\text{pix}$. At
1229 each anterior-posterior (AP) and medial-lateral (ML) coordinate of the 3D atlas, we picked up
1230 the RGB value of the most dorsal brain surface. We obtained a dorsal view of the atlas by
1231 projecting these dorsal RGB values onto a single 2D plane.

1232 ***Selection of injection data***

1233 In the Allen Mouse Brain Connectivity Atlas, each injection experiment is labelled with the name
1234 of the injected structure. First, we narrowed injection experiments using these annotations. We
1235 selected experiments with virus injections into retrosplenial area (RSP), anterior area (VISA) of
1236 posterior parietal association area (PTLp), and secondary motor area (MOs). Then, we further
1237 narrowed down injection experiments based on the exact injection coordinates. As we indicated
1238 in Figure S9, we isolated medial RSP injections, anterior VISA injections, and posterior MOs
1239 injections for RSC, PPC, and pM2, respectively. The database contains experiments that were
1240 performed on wild-type mice and Cre transgenic mice for cell-type specific tracing. We used
1241 experiments from only WT mice or combined data (WT + Cre). The projection patterns were
1242 similar in both cases (Figure S9).

1243 **Axon projection density in dorsal cortex**

1244 We analyzed the axon projection density from RSC, PPC, and pM2 in the dorsal cortex. For
1245 each injection experiment, we calculated the projection density at each AP-ML coordinate as [#
1246 of positive pixels] / [# of all pixels] in the volume of 25 μ m (AP axis) \times 25 μ m (ML axis) \times 1000 μ m
1247 (DV axis, from dorsal surface at each AP-ML coordinate). To create a mean projection density
1248 map, experiments with left hemisphere injections were mirrored relative to midline before
1249 averaging. We also quantified mean projection density within each imaging FOV that we used
1250 for *in vivo* 2-photon calcium imaging. Although our imaging FOVs were based on stereotactic
1251 coordinates from the bregma in the Paxinos' atlas (Paxinos and Franklin, 2004), the Allen
1252 Reference Atlas does not include coordinates from the bregma. To register our imaging FOVs to
1253 the Allen Reference Atlas Coordinate, we calculated the scaling factors for the AP and ML
1254 dimensions of the mouse brain to match the brain in the Paxino's atlas to the brain in the Allen
1255 Reference Atlas. Using the scaling factors, we estimated the coordinates of each imaging FOV
1256 on the Allen Reference Atlas. We calculated the mean signal density within each imaging FOV
1257 of the size 500 μ m \times 500 μ m. The mean signal density of each FOV was used to construct the
1258 connectivity matrix in Figure 8C.

1259

1260 **Data analysis and statistics**

1261 All data analyses and network simulations were performed in Python3.7 with libraries of
1262 TensorFlow (Abadi et al., 2016), scikit-learn (Pedregosa et al., 2011), NumPy (Harris et al.,
1263 2020), SciPy (Virtanen et al., 2020), and Statsmodels (Seabold and Perktold, 2010). Statistical
1264 tests were performed either in Python with SciPy and Statsmodels or in R with its statistics
1265 libraries. All accuracy measures reported in this study were obtained with cross-validation.
1266 Unless otherwise noted, we split trials into training set (70%), validation set (10%), and test set
1267 (20%) for each iteration of decoding, and repeated the network training 5 times by using
1268 different sets of trials as the test set. When we compared ΔQ retrieval/decoding accuracy
1269 across different cortical populations, we matched the number of cells in the input population
1270 activity by subsampling 240 cells in each iteration allowing repetitions with the smallest number
1271 of iterations to include every cell at least once for decoding, and the accuracies from the
1272 iterations were averaged. For all the simulations with artificial population activity, we created 10
1273 distinct populations for each of the 3 types of coding modes by independently sampling tuning
1274 curve slopes from RSC neurons. These repetitions allowed us to tell the variability that
1275 originates from the randomness of ΔQ signal assignments and randomness of network

1276 trainings. All figure plots were created using Matplotlib (Hunter, 2007) and seaborn (Waskom,
1277 2021) in Python.

1278

1279

1280 **References**

1281 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A.,
1282 Dean, J., Devin, M., et al. (2016). TensorFlow: Large-Scale Machine Learning on
1283 Heterogeneous Distributed Systems. 2th USENIX Symp. Oper. Syst. Des. Implement. (OSDI
1284 16), USENIX Assoc. 265–283.

1285 Allen, W.E., Chen, M.Z., Pichamoorthy, N., Tien, R.H., Pachitariu, M., Luo, L., and Deisseroth,
1286 K. (2019). Thirst regulates motivated behavior through modulation of brainwide neural
1287 population dynamics. *Science* (80-). 364.

1288 Banino, A., Barry, C., Uria, B., Blundell, C., Lillicrap, T., Mirowski, P., Pritzel, A., Chadwick, M.J.,
1289 Degrif, T., Moadyil, J., et al. (2018). Vector-based navigation using grid-like representations in
1290 artificial agents. *Nature* 557, 429–433.

1291 Bari, B.A., Grossman, C.D., Lubin, E.E., Rajagopalan, A.E., Cressy, J.I., and Cohen, J.Y.
1292 (2019). Stable Representations of Decision Variables for Flexible Behavior. *Neuron* 103, 922–
1293 933.e7.

1294 Cavanagh, S.E., Towers, J.P., Wallis, J.D., Hunt, L.T., and Kennerley, S.W. (2018). Reconciling
1295 persistent and dynamic hypotheses of working memory coding in prefrontal cortex. *Nat.*
1296 *Commun.* 9.

1297 Chen, T.W., Wardill, T.J., Sun, Y., Pulver, S.R., Renninger, S.L., Baohan, A., Schreiter, E.R.,
1298 Kerr, R.A., Orger, M.B., Jayaraman, V., et al. (2013). Ultrasensitive fluorescent proteins for
1299 imaging neuronal activity. *Nature* 499, 295–300.

1300 Churchland, M.M., Cunningham, J.P., Kaufman, M.T., Foster, J.D., Nuyujukian, P., Ryu, S.I.,
1301 Shenoy, K. V., and Shenoy, K. V. (2012). Neural population dynamics during reaching. *Nature*
1302 487, 51–56.

1303 Evangelidis, G.D., and Psarakis, E.Z. (2008). Parametric image alignment using enhanced
1304 correlation coefficient maximization. *IEEE Trans. Pattern Anal. Mach. Intell.* 30, 1858–1865.

1305 Friedrich, J., Zhou, P., and Paninski, L. (2017). Fast online deconvolution of calcium imaging
1306 data. *PLoS Comput. Biol.* **13**.

1307 Fuster, J.M., and Alexander, G.E. (1971). Neuron activity related to short-term memory. *Science*
1308 (80-). **173**, 652–654.

1309 Guo, Z. V., Inagaki, H.K., Daie, K., Druckmann, S., Gerfen, C.R., and Svoboda, K. (2017).
1310 Maintenance of persistent activity in a frontal thalamocortical loop. *Nature* **545**, 181–186.

1311 Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D.,
1312 Wieser, E., Taylor, J., Berg, S., Smith, N.J., et al. (2020). Array programming with NumPy.
1313 *Nature* **585**, 357–362.

1314 Hattori, R. (2021). PatchWarp. Zenodo. <https://doi.org/10.5281/zenodo.5590965>.

1315 Hattori, R., and Hensch, T.K. (2017). Developmental dynamics of cross-modality in mouse
1316 visual cortex. *BioRxiv* 150847.

1317 Hattori, R., Südhof, T.C., Yamakawa, K., and Hensch, T.K. (2017). Enhanced cross-modal
1318 activation of sensory cortex in mouse models of autism. *BioRxiv* 150839.

1319 Hattori, R., Danskin, B., Babic, Z., Mlynaryk, N., and Komiyama, T. (2019). Area-Specificity and
1320 Plasticity of History-Dependent Value Coding During Learning. *Cell* **177**.

1321 Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.R. (2012).
1322 Improving neural networks by preventing co-adaptation of feature detectors. *ArXiv*.

1323 Hunter, J.D. (2007). Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.* **9**, 90–95.

1324 Inagaki, H.K., Fontolan, L., Romani, S., and Svoboda, K. (2019). Discrete attractor dynamics
1325 underlies persistent activity in the frontal cortex. *Nature* **566**, 212–217.

1326 Iurilli, G., Ghezzi, D., Olcese, U., Lassi, G., Nazzaro, C., Tonini, R., Tucci, V., Benfenati, F., and
1327 Medini, P. (2012). Sound-Driven Synaptic Inhibition in Primary Visual Cortex. *Neuron* **73**, 814–
1328 828.

1329 Jung, Y., Kennedy, A., Chiu, H., Mohammad, F., Claridge-Chang, A., and Anderson, D.J.
1330 (2020). Neurons that Function within an Integrator to Promote a Persistent Behavioral State in
1331 *Drosophila*. *Neuron* **105**, 322-333.e5.

1332 Kennedy, A., Kunwar, P.S., Li, L. yun, Stagkourakis, S., Wagenaar, D.A., and Anderson, D.J.
1333 (2020). Stimulus-specific hypothalamic encoding of a persistent defensive state. *Nature* **586**,

1334 730–734.

1335 Koay, S.A., Thibierge, S., Brody, C.D., and Tank, D.W. (2020). Amplitude modulations of cortical
1336 sensory responses in pulsatile evidence accumulation. *Elife* 9.

1337 Kobak, D., Brendel, W., Constantinidis, C., Feierstein, C.E., Kepcs, A., Mainen, Z.F., Qi, X.L.,
1338 Romo, R., Uchida, N., and Machens, C.K. (2016). Demixed principal component analysis of
1339 neural population data. *Elife* 5.

1340 Kriegeskorte, N. (2015). Deep Neural Networks: A New Framework for Modeling Biological
1341 Vision and Brain Information Processing. *Annu. Rev. Vis. Sci.* 1, 417–446.

1342 Li, N., Daie, K., Svoboda, K., and Druckmann, S. (2016). Robust neuronal dynamics in premotor
1343 cortex during motor planning. *Nature* 532, 459–464.

1344 Lillicrap, T.P., Santoro, A., Marris, L., Akerman, C.J., and Hinton, G. (2020). Backpropagation
1345 and the brain. *Nat. Rev. Neurosci.* 21, 335–346.

1346 Maas, A., Le, Q., O’Neil, T., Vinyals, O., Nguyen, P., and Ng, A. (2012). Recurrent neural
1347 networks for noise reduction in robust ASR. *INTERSPEECH*.

1348 Mante, V., Sussillo, D., Shenoy, K. V., and Newsome, W.T. (2013). Context-dependent
1349 computation by recurrent dynamics in prefrontal cortex. *Nature* 503, 78–84.

1350 Marques, J.C., Li, M., Schaak, D., Robson, D.N., and Li, J.M. (2020). Internal state dynamics
1351 shape brainwide activity and foraging behaviour. *Nature* 577, 239–243.

1352 Masse, N.Y., Yang, G.R., Song, H.F., Wang, X.J., and Freedman, D.J. (2019). Circuit
1353 mechanisms for the maintenance and manipulation of information in working memory. *Nat.*
1354 *Neurosci.* 22, 1159–1167.

1355 McClelland, J., Rumelhart, D., and PDP Research Group (1986). *Parallel Distributed Processing*
1356 (MIT Press).

1357 Miller, E.K., Erickson, C.A., and Desimone, R. (1996). Neural mechanisms of visual working
1358 memory in prefrontal cortex of the macaque. *J. Neurosci.* 16, 5154–5167.

1359 Mitani, A., and Komiya, T. (2018). Real-time processing of two-photon calcium imaging data
1360 including lateral motion artifact correction. *Front. Neuroinform.* 12.

1361 Murray, J.D., Bernacchia, A., Roy, N.A., Constantinidis, C., Romo, R., and Wang, X.J. (2017).
1362 Stable population coding for working memory coexists with heterogeneous neural dynamics in

1363 prefrontal cortex. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 394–399.

1364 Musall, S., Kaufman, M.T., Juavinett, A.L., Gluf, S., and Churchland, A.K. (2019). Single-trial
1365 neural dynamics are dominated by richly varied movements. *Nat. Neurosci.* **22**, 1677–1686.

1366 Oh, S.W., Harris, J.A., Ng, L., Winslow, B., Cain, N., Mihalas, S., Wang, Q., Lau, C., Kuan, L.,
1367 Henry, A.M., et al. (2014). A mesoscale connectome of the mouse brain. *Nat.* **2014** 5087495
1368 **508**, 207–214.

1369 Orhan, A.E., and Ma, W.J. (2019). A diverse range of factors affect the nature of neural
1370 representations underlying short-term memory. *Nat. Neurosci.* **22**, 275–283.

1371 Osa, T., Pajarinen, J., Neumann, G., Bagnell, J.A., Abbeel, P., and Peters, J. (2018). An
1372 Algorithmic Perspective on Imitation Learning. *Found. Trends Robot.* **7**, 1–179.

1373 Pachitariu, M., Stringer, C., Dipoppa, M., Schröder, S., Rossi, L.F., Dalgleish, H., Carandini, M.,
1374 and Harris, K. (2016). Suite2p: beyond 10,000 neurons with standard two-photon microscopy.
1375 *BioRxiv* 061507.

1376 Pachitariu, M., Stringer, C., and Harris, K.D. (2018). Robustness of spike deconvolution for
1377 neuronal calcium imaging. *J. Neurosci.* **38**, 7976–7985.

1378 Pascanu, R., Mikolov, T., and Bengio, Y. (2012). On the difficulty of training Recurrent Neural
1379 Networks. *30th Int. Conf. Mach. Learn. ICML 2013* 2347–2355.

1380 Paxinos, G., and Franklin, K.B.J. (2004). *The Mouse Brain in Stereotaxic Coordinates*, 2nd
1381 edition. Acad. Press 360 p.

1382 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
1383 Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine Learning in Python.
1384 *J. Mach. Learn. Res.* **12**, 2825–2830.

1385 Rogers, T.T., and McClelland, J.L. (2014). Parallel distributed processing at 25: Further
1386 explorations in the microstructure of cognition. *Cogn. Sci.* **38**, 1024–1077.

1387 Romo, R., Brody, C.D., Hernández, A., and Lemus, L. (1999). Neuronal correlates of parametric
1388 working memory in the prefrontal cortex. *Nature* **399**, 470–473.

1389 Rumelhart, D., McClelland, J., and PDP Research Group (1986). *Parallel Distributed Processing*
1390 (MIT Press).

1391 Russo, A.A., Bittner, S.R., Perkins, S.M., Seely, J.S., London, B.M., Lara, A.H., Miri, A.,

1392 Marshall, N.J., Kohn, A., Jessell, T.M., et al. (2018). Motor Cortex Embeds Muscle-like
1393 Commands in an Untangled Population Response. *Neuron* **97**, 953–966.e8.

1394 Russo, A.A., Khajeh, R., Bittner, S.R., Perkins, S.M., Cunningham, J.P., Abbott, L.F., and
1395 Churchland, M.M. (2020). Neural Trajectories in the Supplementary Motor Area and Motor
1396 Cortex Exhibit Distinct Geometries, Compatible with Different Classes of Computation. *Neuron*
1397 **107**, 745–758.e6.

1398 Seabold, S., and Perktold, J. (2010). Statsmodels: Econometric and Statistical Modeling with
1399 Python. *PROC. 9th PYTHON Sci. CONF.*

1400 Serences, J.T. (2008). Value-Based Modulations in Human Visual Cortex. *Neuron* **60**, 1169–
1401 1181.

1402 Song, H.F., Yang, G.R., and Wang, X.J. (2017). Reward-based training of recurrent neural
1403 networks for cognitive and value-based tasks. *Elife* **6**.

1404 Srivastava, N., Hinton, G., Krizhevsky, A., and Salakhutdinov, R. (2014). Dropout: A Simple
1405 Way to Prevent Neural Networks from Overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958.

1406 Steinmetz, N.A., Zatka-Haas, P., Carandini, M., and Harris, K.D. (2019). Distributed coding of
1407 choice, action and engagement across the mouse brain. *Nature* **576**, 266–273.

1408 Stringer, C., Pachitariu, M., Steinmetz, N., Reddy, C.B., Carandini, M., and Harris, K.D. (2019).
1409 Spontaneous behaviors drive multidimensional, brainwide activity. *Science* (80-.). **364**.

1410 Sutton, R.S., and Barto, A.G. (2018). Reinforcement Learning: An Introduction (2nd ed.) (MIT
1411 Press).

1412 Tsuda, B., Tye, K.M., Siegelmann, H.T., and Sejnowski, T.J. (2020). A modeling framework for
1413 adaptive lifelong learning with transfer and savings through gating in the prefrontal cortex. *Proc.*
1414 *Natl. Acad. Sci. U. S. A.* **117**, 29872–29882.

1415 Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P.-A. (2010). Stacked
1416 Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local
1417 Denoising Criterion. *J. Mach. Learn. Res.* **11**, 3371–3408.

1418 Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski,
1419 E., Peterson, P., Weckesser, W., Bright, J., et al. (2020). SciPy 1.0: fundamental algorithms for
1420 scientific computing in Python. *Nat. Methods* **17**, 261–272.

1421 Wang, J.X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J.Z., Hassabis, D.,
1422 and Botvinick, M. (2018). Prefrontal cortex as a meta-reinforcement learning system. *Nat.*
1423 *Neurosci.* **21**, 860–868.

1424 Waskom, M.L. (2021). seaborn: statistical data visualization. *J. Open Source Softw.* **6**, 3021.

1425 Wekselblatt, J.B., Flister, E.D., Piscopo, D.M., and Niell, C.M. (2016). Large-scale imaging of
1426 cortical dynamics during sensory perception and behavior. *J. Neurophysiol.* **115**, 2852–2866.

1427 Yamins, D.L.K., and DiCarlo, J.J. (2016). Using goal-driven deep learning models to understand
1428 sensory cortex. *Nat. Neurosci.* **2016** **19**, 356–365.

1429 Zhu, J., Cheng, Q., Chen, Y., Fan, H., Han, Z., Hou, R., Chen, Z., and Li, C.T. (2020). Transient
1430 Delay-Period Activity of Agranular Insular Cortex Controls Working Memory Maintenance in
1431 Learning Novel Tasks. *Neuron* **105**, 934-946.e5.

1432