

1 **FRONT MATTER**

2 **Title**

- 3 • Beyond the Delay Neural Dynamics: a Decoding Strategy for Working Memory Error
4 Reduction

5 **Authors**

6 Zeyuan Ye^{1,2,3,4,9}, Haoran Li^{1,9}, Liang Tian^{1,5,6,7*}, Changsong Zhou^{1,2,5,8,10*}

7 **Affiliations**

8 ¹Department of Physics, Hong Kong Baptist University, Hong Kong, China

9 ²Centre for Nonlinear Studies and Beijing-Hong Kong-Singapore Joint Centre for Nonlinear
10 and Complex Systems (Hong Kong), Hong Kong Baptist University, Hong Kong, China

11 ³Institute of Interdisciplinary Studies, Hong Kong Baptist University, Hong Kong, China

12 ⁴Department of Physics, Washington University in St. Louis, St. Louis, Missouri, USA

13 ⁵Institute of Computational and Theoretical Studies, Hong Kong Baptist University, Hong
14 Kong, China

15 ⁶Institute of Systems Medicine and Health Sciences, Hong Kong Baptist University, Hong
16 Kong, China

17 ⁷State Key Laboratory of Environmental and Biological Analysis, Hong Kong Baptist
18 University, Hong Kong, China

19 ⁸Life Science Imaging Centre, Hong Kong Baptist University, Hong Kong, China

20 ⁹These authors contributed equally.

21 ¹⁰Lead contact

22 *Correspondence and requests for materials should be addressed to L.T. or C.Z. (email:
23 liangtian@hkbu.edu.hk; cszhou@hkbu.edu.hk)

24 **Abstract**

25 Understanding how the brain preserves information despite intrinsic noise is a fundamental
26 question in working memory. Typical working memory tasks consist of delay phase for
27 maintaining information, and decoding phase for retrieving information. While previous
28 works have focused on the delay neural dynamics, it is poorly understood whether and how
29 the neural process during decoding phase reduces memory error. We studied this question
30 by training recurrent neural networks (RNNs) on a color delayed-response task. We found
31 that the trained RNNs reduce the memory error of high-probability-occurring colors
32 (common colors) by decoding/attributing a broader range of neural states to them during
33 decoding phase. This decoding strategy can be further explained by a continuing converging
34 neural dynamics following delay phase and a non-dynamic biased readout process. Our
35 findings highlight the role of the decoding phase in working memory, suggesting that neural
36 systems deploy multiple strategies across different phases to reduce memory errors.

42 **Significance**

43 Preserving information under noise is crucial in working memory. A typical working memory
44 task consists of a delay phase for maintaining information, and a decoding phase for decoding the
45 maintained into an output action. While the delay neural dynamics have been intensively studied,
46 the impact of the decoding phase on memory error reduction remains unexplored. We trained
47 recurrent neural networks (RNNs) on a color delayed-response task and found that RNNs reduce
48 memory error of a color by decoding a larger portion of the neural state to that color. This strategy
49 is supported both by a converging neural dynamic, and a non-dynamic readout process. Our
50 results suggest that neural networks can utilize diverse strategies, beyond delay neural dynamics,
51 to reduce memory errors.

52

53 **MAIN TEXT**

54

55 **Introduction**

56

57 Working memory is the ability to maintain information for a short period of time without external
58 stimuli. It largely consists of a perception phase for sensing the information, a delay phase for
59 maintaining information, and finally, a decoding phase (for example, response epoch in typical
60 delayed-response tasks¹⁻³) for retrieving information⁴. Working memory tasks are fundamentally
61 challenging because the neural system needs to maintain accurate information despite intrinsic
62 stochastic noise. Without additional error-correcting mechanisms, the neural population activity
63 will deviate from its original state, leading to large memory errors⁵⁻⁸. The neural mechanisms
64 utilized by the neural system to mitigate these memory errors remain unclear.

65

66 Previous works have primarily focused on the neural dynamics during the delay phase^{1,5-14}. For
67 instance, the neural system can form discretized attractors to represent some information values^{1,9-}
68 ¹¹. An attractor is a special neural population state; any small deviation from the attractor state will
69 be brought back. This unique property of attractors can stabilize neural population states against
70 random deviations due to noise, thereby reducing memory errors. This attractor-based memory
71 error correcting mechanism has been supported by experimental behavior data¹, experimental
72 neural recordings¹⁵, and simulations from artificial neural networks^{7,11,16}. It can also be
73 implemented by artificial neural networks using biologically plausible synaptic rules^{9,17}.

74

75 However, despite also being a key phase, the role of the decoding phase has largely been neglected.
76 From the information processing perspective, the decoding phase acts as a decoding mapping,
77 which maps the maintained delayed neural state to an output action¹⁸⁻²³. A change in the decoding
78 mapping will lead to a significant change of behavioral performance. This importance of decoding
79 mapping can be demonstrated in the Brain-Computer Interface (BCI) experiments^{23,24}, where the
80 BCI functions as a decoder, mapping neural population states to external actions (e.g. cursor
81 movement). Alterations in the BCI's decoding mapping would lead to significant action errors,
82 necessitating relearning of the new decoding pattern by the animal^{23,24}. Similarly, in working
83 memory, the delay neural population state must progress through decoding phase, be decoded to
84 muscles signals for an output action (e.g. saccade). What is the decoding mapping from an end-of-
85 -delay neural state to an output information? Whether and how such mapping helps reducing the
86 memory error? What processes occur during the decoding phase to establish this decoding
87 mapping? Addressing these questions is important for understanding the diverse strategies, beyond
88 just delay neural dynamics, of the neural system to reduce the memory error.

89

90 In this paper, we trained artificial recurrent neural networks (RNNs) to perform a color delayed-
91 response task, where the color in each trial was sampled from a prior distribution with a few high-

92 probability colors (common colors)^{11,25,26}. We found that the trained RNN exhibited smaller
93 memory error on common colors, which aligns with previous behavioral experiments¹. We found
94 two main mechanisms that the RNNs used to reduce memory error. First, the neural system created
95 attractors to encode common colors during the delay phase. Second, during the decoding phase, a
96 large part of the neural population states was decoded to common colors, which improved noise
97 tolerance. This noise-tolerant decoding mapping can be further understood by (1) continuing the
98 attractor-based dynamics during the decoding phase, and (2) a non-dynamic biased readout from
99 the recurrent neural state to common colors. Further, we proposed an approximation formula which
100 naturally decompose the memory error into delay dynamic and decoding components – neglecting
101 decoding component will lead to a failure of explaining the RNN’s memory error. Our results
102 emphasize the importance of the decoding phase and propose experimentally testable predictions.

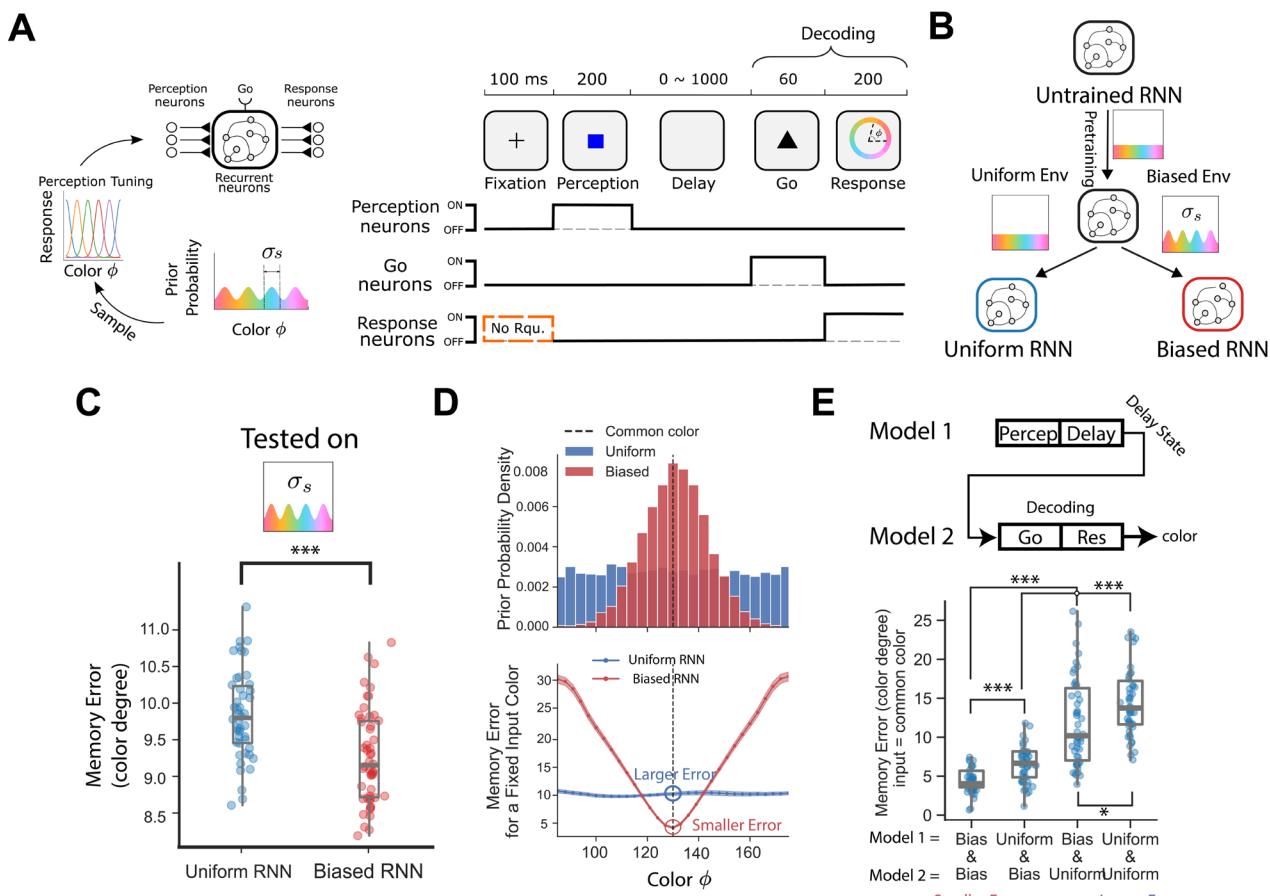
103 **Results**

104 **RNNs trained in a biased environment reproduce human behavioural features in the color
105 delayed-response task.**

106 Our study utilized a vanilla RNN architecture, which included 12 perception neurons, one go
107 neuron, 256 recurrent neurons, and 12 response neurons (Figure 1A and Methods). We added noise
108 to the recurrent neurons to simulate intrinsic brain noise. Task consists of a fixation, a perception,
109 a delay, a go cue and finally a response epoch¹. Specifically, in each trial, we sampled a color from
110 an environmental prior distribution with four peaks corresponding to common colors (similar to the
111 behavioral experiment¹). The sampled input color was perceived by the perception neurons through
112 predefined tuning curves (von Mises functions with shifting mean positions, see Figure 1A), and
113 we added Gaussian noise perturbations to simulate perception noise (see Methods). The perception
114 neurons then fed the activity forward to the 256 recurrent neurons. During the delay period
115 (randomly determined delay epoch duration), the RNN used these recurrent neurons to maintain the
116 color information. After the delay, the go neuron received a signal that initiated the response action.
117 The response neurons, equal in number to the perception neurons, activated during the response
118 epoch, were expected to mirror the activity of the perception neurons during the perception epoch
119 (see Methods). We then averaged the activity of these response neurons over time and mapped it
120 back into color information using the population vector method. Go and response epochs in
121 combination is considered as a decoding phase in this paper.

122
123 Each RNN was trained progressively (Figure 1B and Method)²⁷. The progressive training consisted
124 of two stages. First, in the pretraining stage, the randomly initialized RNN was trained in a uniform-
125 prior environment. Following this, the pretrained RNN was retrained in a new environment. If the
126 new environmental prior is biased, the retrained RNN is referred to as the "Biased RNN".
127 Conversely, if the new environmental prior is uniform, the trained RNN is called the "Uniform
128 RNN". Our procedure ensured that both types of RNNs had the same number of training trials.

129
130 Similar training with biased environments has been performed experimentally before (with
131 humans)^{1,28}. We then asked whether our Biased RNN could reproduce some basic behavioral
132 features of those experiments. We found that (see SI Figure 1), in alignment with previous
133 experiments¹, (1) The RNN's memory error increases with longer delays; (2) Despite using
134 uniformly sampled colors as inputs, Biased RNN still outputs biased color distribution, with four
135 output color distribution peaks aligned with the four biased environmental priors in the training; (3)
136 When the input trial's color is near the common color, the output color of the Biased RNN tends to
137 shift closer to the common colors, suggesting that Biased RNNs have a tendency to output common
138 colors.



140
141 **Fig 1. Decoding phase (go and response epochs) is crucial for RNNs to reduce memory error.** (A) An
142 environment was defined as a probability prior function for the input color. In each trial, a color was sampled
143 from the prior (four common colors = $40^\circ, 130^\circ, 220^\circ, 310^\circ$), then sensed by the perception neurons.
144 Following the perception epoch, the RNN went through the delay epoch for maintain information (delay
145 length was sampled randomly from 0 to 1000 ms), go epoch to be ready, and finally response epoch to
146 reproduce the previously sensed color. Except regularization terms, loss function was only applied on the
147 response neurons, which had no requirement in fixation epoch (No Rqu.) and should be silent in the
148 perception, delay and go epochs. (B) An RNN was trained by multiple trials progressively (see texts and
149 Methods). Depending on the retrained environment prior, the final trained RNN was categorized as a
150 Uniform RNN (uniform prior) or a Biased RNN (biased prior). (C) After training, RNNs were evaluated
151 under a biased environment ($\sigma_s = 12.5^\circ$, same as the training environment of the Biased RNN). The memory
152 error of each RNN was calculated as the root-mean-squared difference between output and input colors over
153 5000 trials. Each dot represents one RNN, with a total of 50 RNNs used for each type. ($***: p < 10^{-3}$
154 Wilcoxon signed-rank test). Outlier RNNs were not shown. (D) Top: The prior distribution of colors around
155 a common color (130° , black dashed line). Bottom: Memory error for various input colors, assessed over
156 5000 trials for each input color and RNN. The colored lines represent the mean memory errors across the 50
157 RNNs, with error bands indicating standard errors. Outlier RNNs were removed (see Methods). (E) Cross
158 decoding suggests the importance of decoding phase in reducing memory error (e.g., Bias & Bias vs. Bias
159 & Uniform). Delay recurrent neural state was prepared from a model 1, then was decoded by model 2. Each
160 dot is one randomly sampled model1-model2 pair. 50 pairs were sampled for each model1-model2 category
161 combination. $***: p < 10^{-3}; *: p < 5 \times 10^{-2}$; Wilcoxon signed-rank test. Boxes indicate the interquartile
162 range between the first and third quartiles with the central mark inside each box indicating the median.
163 Whiskers extend to the lowest and highest values within 1.5 times the interquartile range. Outlier RNNs
164 were not shown. Biased RNN was trained under $\sigma_s = 12.5^\circ$ in panel (C, D, E).

167 **Cross-decoding experiments showed that both delay and decoding phases are crucial in**

168 **memory error reduction.**

169 Using the trained Biased and Uniform RNNs, we then explored the properties of memory error.
170 Memory error of a trial is defined as the circular subtraction (circular period is 360°) of the RNN's
171 output color from the trial's input color. The memory error across multiple trials is defined as the
172 root-mean-square error of individual trial errors.

173 First, we evaluated the RNN's overall memory error in a biased environment ($\sigma_s = 12.5^\circ$) which
174 is the same environment the Biased RNNs were trained in. As expected, the Biased RNN exhibited
175 smaller memory errors (Figure 1C, $p < 10^{-3}$, Wilcoxon signed-rank test). Mathematically in the
176 limit of large number of trials, the memory error is the weighted average (weighted by the prior) of
177 the memory error for every fixed input color. Hence, we examined the memory error for each fixed
178 input color (Figure 1D). For the Uniform RNN, memory errors were uniform across all input colors.
179 In contrast, the Biased RNN showed a reduced memory error for the common color (e.g., at 130°)
180 compared to other input colors and also to the same input color in the Uniform RNN. This reduction
181 of memory error on common colors provided an opportunity to explore the RNN's neural
182 mechanisms underlying memory error reduction.

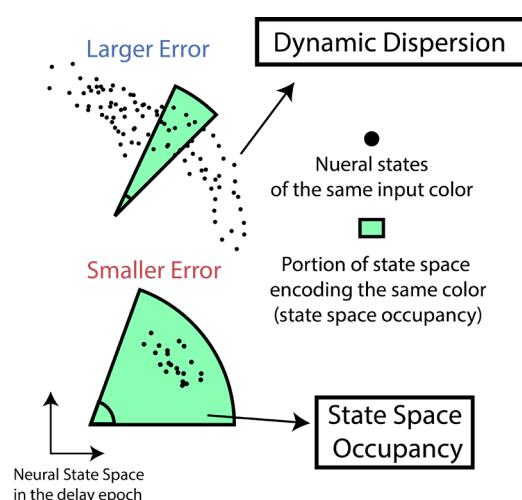
183 Previous studies have primarily focused on delay dynamic mechanisms to reduce memory errors¹.
184 We hypothesized that, in addition to neural dynamics during the delay, the decoding from the end-
185 of-delay neural state to an output color during decoding phase is also crucial. To test this, we
186 conducted a cross-decoding experiment, where the non-decoding (fixation, perception, and delay
187 epochs) and decoding phases (go and response epochs) were separated (see Figure 1E, upper panel).
188 In this experiment, two RNN models were randomly chosen from either the Biased or Uniform
189 category (two models can belong to the same or different categories). Model 1 underwent a trial
190 with a fixed common color input, and the recurrent neural state (i.e. recurrent neural population
191 activity, or neural state for short) at the delay's end was collected. Collected recurrent neural state
192 was then used to set the recurrent neural state of model 2. This step requires a one-to-one matching
193 from model 1's recurrent neurons to model 2's recurrent neurons. This matching was achieved by
194 comparing the preferred colors of the individual neurons between the two neural populations (see
195 Methods). Model 2 then directly ran through the decoding phase (go and response epochs) to output
196 a color. The error between the output and input color across 500 trials was the memory error of this
197 single RNN pairing (model 1-model 2). This entire process was conducted for 50 random pairs for
198 each category combination (Biased & Uniform, Biased & Biased, Uniform & Uniform, Uniform &
199 Biased). It should be noted that the algorithm used for matching (comparing preferred colors) from
200 the recurrent neural state of model 1 to model 2 is not optimal because the connections are not
201 exactly the same between RNN models. However, we applied the same matching algorithm to all
202 category combinations, so comparisons between different category combinations should be fair to
203 a certain extent.

204 We found that (Figure 1E), for the same type of delay state preparation RNNs (model 1), decoding
205 using Biased RNNs significantly reduced memory error of common colors compared to decoding
206 using Uniform RNNs (i.e., Bias & Bias vs. Bias & Uniform, and Uniform & Bias vs. Uniform &
207 Uniform). This suggests that Biased RNNs learned certain neural mechanisms during the decoding
208 phase to reduce the memory error of common colors.

213 Besides, we found that for a fixed type of decoding RNN (model 2), using Biased RNN for delay
214 state preparation also results in smaller memory errors of common colors. This finding supports the
215 idea that Biased RNN learned neural mechanisms during the delay phase to reduce memory errors.

216 **A hypothesis of the memory error reduction by both delay neural dynamics and decoding**
217 **strategy**

218 The cross-decoding experiments conducted underscore the pivotal roles of both the delay and
219 decoding phases (go cue and response epochs) in influencing memory errors. Previous hypotheses
220 have suggested that reduced neural dispersion dynamics during the delay might contribute to error
221 diminution. In this paper, we introduced a novel hypothesis concerning the importance of the
222 decoding phase in error reduction. Conceptually, the state of a neural population can be represented
223 as a point within a multidimensional state space, where each axis represents the activity of an
224 individual neuron (Figure 2). The end-to-end effect of decoding phase is to transform the neural
225 state at the end of the delay into an output color. Geometrically, decoding phase “maps” distinct
226 portions in the state space to specific output colors. Imagine a scenario where neural states, despite
227 being noisy and dispersed, fall within a region of the state space mapped to a single color; the output
228 color remains constant despite the dispersion of neural states. Thus, colors with larger state space
229 occupancies exhibit greater tolerance to noise, subsequently diminishing memory errors. In
230 summary, we propose that memory error reduction results from two factors (Figure 2): (1) reducing
231 the dynamic dispersion of neural population states during the delay, and (2) assigning a broader
232 range of neural states to common colors during the decoding phase, a process we term 'larger state
233 space occupancy by a color'. This paper primarily tests and explores the second mechanism,
234 focusing on the effect of decoding phase on memory error reduction.
235
236

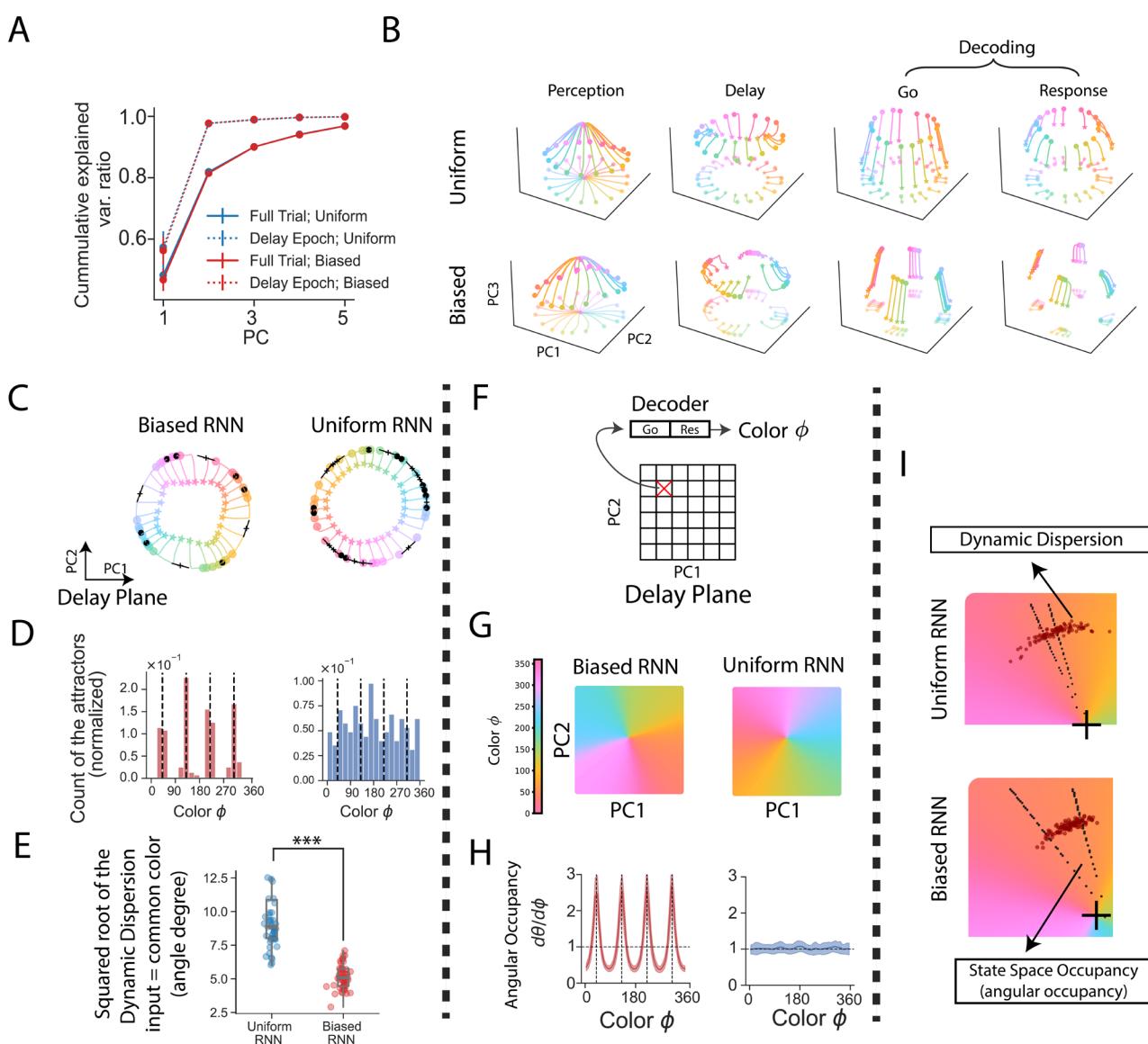


237
238 **Figure 2. A hypothesis decomposing the memory error reduction mechanism into delay neural**
239 **dynamics (dynamic dispersion) and a decoding strategy (state space occupancy).** Illustration of the
240 hypothesis: (1) Each axis is one recurrent neural activity. Each dot is one recurrent neural population state
241 at the end of the delay under a fixed input color. Memory error can be reduced if the neural states are more
242 stable (smaller dynamic dispersion). (2) Green region represents portion of state space that will be decoded
243 into a color. Larger state space occupancy of one color allows better tolerance of neural state dispersion,
244 hence smaller memory error for that color.

245 **RNN's neural activity is low-dimensional.**

246 To test our hypothesis on memory error reduction, we inspected the neural activities in the RNNs.
 247 We conducted multiple trials for each RNN (intrinsic noise was turned off), using uniformly
 248 sampled color inputs. Recurrent neural activities were collected and projected onto the first few
 249 principal components (PCs) using Principal Component Analysis (PCA). The cumulative variance
 250 explained ratio suggests that the neural population activity is essentially low-dimensional (Figure
 251 3A): requiring 3 dimensions to account for over 90% of the variance in neural activity across the
 252 entire trial, and only 2 dimensions when focusing exclusively on the delay epoch.
 253

254 The inherent low dimensionality of the neural activity allowed for direct visualization by projecting
 255 the activity onto the first three PCs, as shown in Figure 3B. During the perception epoch, the neural
 256 states are driven by the input colors, moving towards different directions within a ring-shaped
 257 manifold. In the subsequent delay phase, these neural states remain relatively stable, although slight
 258 lateral movements are observed. Next, in the go epoch, the neural state shifts along the third
 259 principal component. Finally, in the response epoch, the neural states roughly stay stable in a 2D
 260 ring-like structure. These observations from the various phases of neural activity illustrate the
 261 dynamic yet structured nature of neural state dynamics in different phases of task execution.
 262
 263



265 **Figure 3. Reverse engineering of the RNNs shows evidence for the hypothesized memory reduction**
266 **mechanisms (smaller dynamic dispersion and larger state space occupancy).** Biased RNN was trained
267 under $\sigma_s = 12.5^\circ$. **(A)** Dimensionality of the neural activity (1000 trials for each RNN). Dots and error bar
268 represent mean and standard deviation across 50 RNNs. Uniform: Uniform RNN; Biased: Biased RNN. Blue
269 and red traces mostly overlap. **(B)** PCA visualization of neural states. Stars denote the beginning of the
270 epochs; dots represent the end. Each trajectory is one trial with color indicating the trial input color. Both
271 3D trajectories (upper) and their 2D projections (lower) were shown in each panel. **(C)** 2D visualization of
272 two example RNN's neural states in the delay epoch. Black dots are the attractor points. Crosses are saddle
273 points whose long bars indicating the unstable (positive eigenvalue) directions. **(D)** Distribution of attractors.
274 More attractors were formed near common colors (dash lines) in the trained Biased RNN. Results
275 concatenated attractors from 50 RNNs. **(E)** End-of-delay neural states have smaller dynamic dispersion (see
276 Methods) in Biased RNN when the input color is fixed to common color. Each dot is one RNN. ***: $p <$
277 0.001 Wilcoxon signed-rank test. Boxes indicate the interquartile range between the first and third quartiles
278 with the central mark inside each box indicating the median. Whiskers extend to the lowest and highest
279 values within 1.5 times the interquartile range. Outlier RNNs were not shown. **(F)** Decoding the delay plane
280 by selecting mesh points and running through go and response epochs. **(G)** Decoding results of two example
281 RNNs. Color is the decoded color. **(H)** Angular occupancy--the size of angle (θ) on the delay plane used for
282 encoding one unit of color. Solid lines are the mean across 50 RNNs, and error band is the standard deviation.
283 Outlier RNNs were removed (see Methods). **(I)** Two example RNNs illustrate the joint effect of dynamic
284 dispersion and angular occupancy. Each black dot is the position of neural state (on the delay PC1-PC2
285 plane) at the end of delay of one trial (delay length = 800 ms), where trial's input color was fixed to a
286 common color (40°). Dashed lines enclose the angular space for encoding 35° to 45°. Black crosses are the
287 center of the PC1-PC2 plane. In this figure, RNN's intrinsic noise was turned off for better visualization,
288 except in (E, I).

289 **Biased RNN formed attractors near common colors to reduce the dynamic dispersion.**

290 Next, we tested whether the dynamic dispersion during the delay phase is a key factor in memory
291 error reduction. We projected the delay neural activity to the first two PCs. We refer this 2D
292 representation as the “delay plane”. On the delay plane, neural states for different input colors
293 formed a ring-like structure (Figure 3C).

294 We observed small lateral motion of neural states along the ring. According to the dynamic theory,
295 this small motion may be driven by fixed points on the ring. To test this, we searched fixed points
296 by finding the local minimum of speed (see Methods). We found two types of fixed points on the
297 ring: (1) Attractors which attract nearby states (2) Saddle points, which repel nearby neural states
298 along the direction of positive eigenvalues. Interestingly, we found in the Biased RNN, attractors
299 predominantly appeared in four positions (Figure 3C). These attractors also had more negative
300 eigenvalues (indicating stronger attractive force) in more Biased RNNs (with smaller prior σ_s , SI
301 Figure 2).

303 The four dominant positions of four attractors may represent the four common colors in the biased
304 environmental prior. To test this, we developed an RNN decoder. The RNN decoder itself is a
305 faithful copy of the original RNN. Given a neural state of interest, the RNN decoder's recurrent
306 state is set to match the neural state, then without any delay, it directly runs through the go and
307 response epoch. The output color is the decoded color of the neural state of interest.

309 We decoded the attractors of Uniform RNNs and Biased RNNs. The results revealed that in Biased
310 RNNs, attractors were primarily mapped to common colors (Figure 3D). Since attractors have
311 “attracting effects,” this suggests that the neural state may have a smaller dynamic dispersion when
312 the input color is a common color. To test this, we directly measured dynamic dispersion as follows:
313 We conducted multiple trials for each RNN, keeping the input color fixed at the common color and

315 the fixed the delay length at 800 ms. Neural states at the end of the delay were collected, and their
316 angles on the delay plane were computed. The variance of these neural state angles represents the
317 dynamic dispersion. The results (Figure 3E) indicate that the dynamic dispersion for the common
318 color in Biased RNNs is smaller than in Uniform RNNs, supporting the hypothesis that RNNs
319 reduce memory error by diminishing dynamic dispersion¹⁰.

320 **Biased RNN allocates larger angular occupancies to common colors**

321 Next, we tested whether the RNNs formed larger state space occupancy to common colors. The
322 end-to-end effect of decoding phase is transforming neural state at the end of delay epoch to an
323 output color. This end-to-end effect is entirely same as the RNN decoder (Figure 3F), which
324 decodes a neural state by continuing to run the neural state through go and response epochs. We
325 used the RNN decoder to decode the delay plane (not the neural state from any particular trials).
326 Example RNNs' results suggested that color information is represented as angles in the delay plane
327 (Figure 3G).

328 We investigated the quantitative relationship between angle and color. We sampled and decoded
329 dense neural states along a ring (see Methods). This provided a numerical relation between the
330 angle to the decoded color. The numerical differentiation of the angle by the color is called angular
331 occupancy. Angular occupancy measures the size of angular space mapped to a unit change of
332 color. Results indicate that, in the Biased RNN, common color has larger angular occupancy (Figure
333 3H). As a baseline comparison, the angular occupancy of the Uniform RNN is uniform. These
334 results suggested that Biased RNN reduces the memory error (for common color) by allocating
335 larger angular space.

336 We used two example RNNs to provide an intuition of how the dynamics and decoding strategy
337 jointly contribute to the reduction of memory error (Figure 3I). We ran each of the example RNNs
338 through 1000 trials with a fixed common color as input (delay length fixed at 800 ms). Neural states
339 at the end of the delay were collected and visualized on the delay plane. Two dashed lines enclosed
340 the portion of space encoding color 35° to 45°, around a common color 40°. It can be observed
341 that, compared to the Uniform RNN, the Biased RNN has a larger angular occupancy and smaller
342 dynamic dispersion, supporting the hypothetical picture in Figure 2.

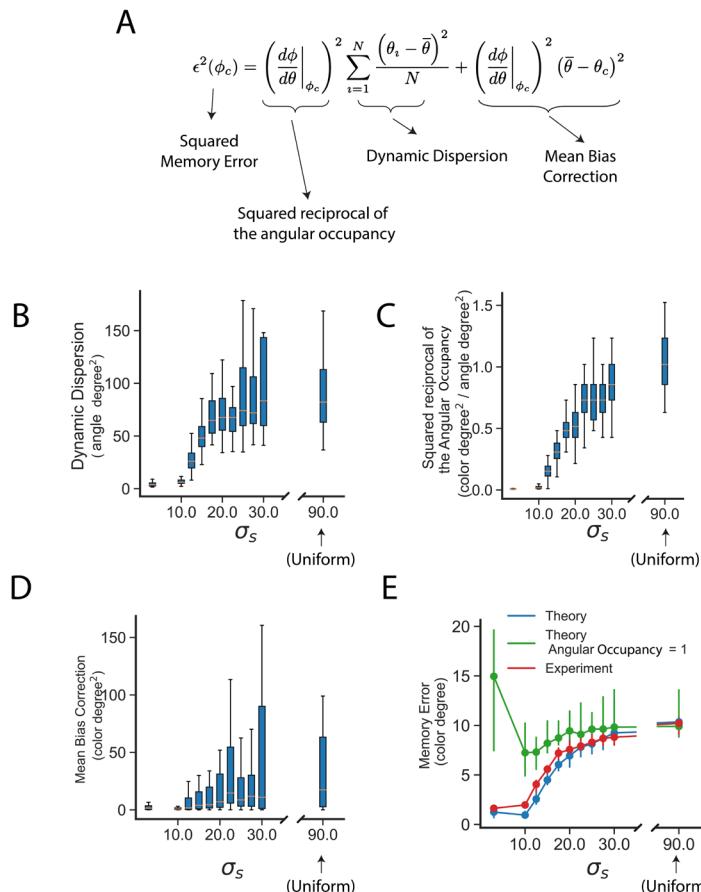
343 **An approximate theory quantitatively relates the dynamic dispersion and angular
344 occupancy to the memory error.**

345 Figure 3I provides a visualization of how the dynamic and decoding strategy jointly contribute to
346 the reduction of memory errors on common colors. Here we study a mathematical conceptualization
347 of how dynamic dispersion and angular occupancy jointly contribute to the reduction of memory
348 errors.

349 The central idea is the Taylor expansion (see Methods). Color ϕ is a function of angle θ in PC1-
350 PC2 space. Denoting common color as ϕ_c , the corresponding angle as θ_c . Assuming that, when
351 trial input is the common color, the actual angle of neural state θ does not deviate from the true
352 color angle θ_c too much. Therefore, we expanded the color function as $\phi(\theta) \approx \phi_c +$
353 $(d\phi/d\theta|_{\phi_c})(\theta - \theta_c)$, where $d\phi/d\theta|_{\phi_c}$ is the reciprocal of angular occupancy. Inserting this
354 approximation into the squared memory error $\epsilon^2(\phi_c) = \sum_{i=1}^N (\phi_i - \phi_c)^2 / N$ (where ϕ_i is the output
355 color in trial i), we can then obtain an approximation formula (Equation 1 and Figure 4A)
356 describing the memory error using angular occupancy ($d\phi/d\theta|_{\phi_c}$), dynamic dispersion
357 ($\sum_{i=1}^N (\theta_i - \bar{\theta})^2 / N$) and a mean bias correction ($(\bar{\theta} - \theta_c)^2$)

$$\epsilon^2(\phi_c) = \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 \sum_{i=1}^N \frac{(\theta_i - \bar{\theta})^2}{N} + \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 (\bar{\theta} - \theta_c)^2, \quad (1)$$

361 where the mean bias correction described the correction due to the shift of trial averaged mean
 362 angle $\bar{\theta}$ to θ_c .



363
 364 **Figure 4. An approximation formula decomposes the memory error of a common color into a dynamic**
 365 **dispersion, angular occupancy, and a correction term. (A)** An approximation formula for the squared
 366 **memory error (see text). (B, C, D)** Dynamic dispersions, squared reciprocal of the angular occupancy and
 367 **mean bias correction for biased RNNs trained on different biased environment (50 RNNs for each σ_s).**
 368 Smaller σ_s indicating narrower prior (see Figure 1A). Boxes indicate the interquartile range between the first
 369 and third quartiles cross 50 RNNs with the central mark inside each box indicating the median. Whiskers
 370 extend to the lowest and highest values within 1.5 times the interquartile range. **(E)** Comparing the theory
 371 to the experimental RNN's memory error. Dots are the median memory error, and error bars show first and
 372 the third quantiles across 50 RNNs for each σ_s . Theoretical prediction was computed as (A). Experimental
 373 memory error of each RNN was measured by computing the averaged memory error of 5000 trials (fixing
 374 input at common color, outlier trial errors were removed). As a comparison, we also computed the prediction
 375 of theory but assuming a trivial angular occupancy (equals to 1).

376
 377 To test this approximation formula (Figure 4A), we trained biased RNNs on different environmental
 378 priors by changing the value of σ_s to control the prior width. The smaller the σ_s , the narrower the
 379 prior, and the higher the prior probability of sampling a common color. We then computed dynamic
 380 dispersion, squared reciprocal of the angular occupancy and mean biased correction separately (see
 381 Methods). On the other hand, we also computed the experimental memory error directly by running
 382 each RNN 5000 trials, fixing input color as common color.

384 We found the higher common color prior probability is (smaller σ_s), the smaller dynamic dispersion
385 and the larger angular occupancy (Figure 4B, C) are. In general, the theoretical prediction has good
386 alignment with the actual experimental memory error (Figure 4E). To test whether the angular
387 occupancy term is important for the memory error, we computed theoretical prediction (Figure 4A)
388 but setting angular occupancy to 1. We showed that neglecting angular occupancy factor will lead
389 to a failure of explaining the RNN's memory error for small σ_s (Figure 4E).

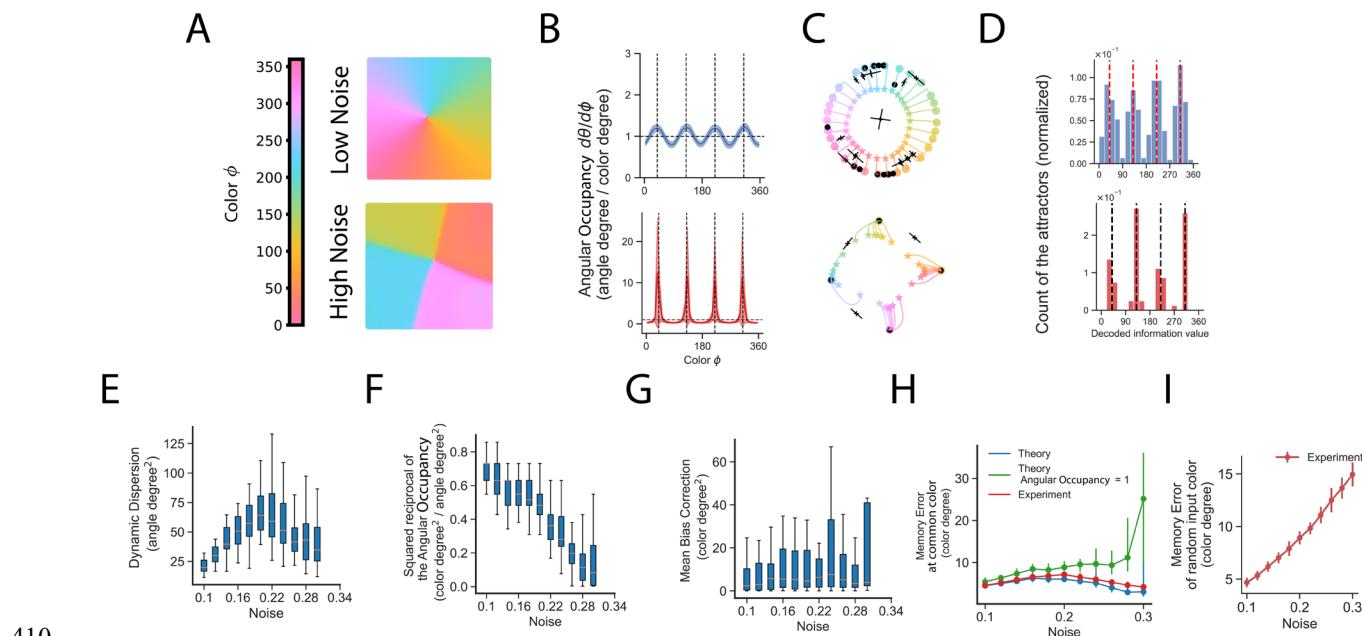
390 **Higher intrinsic noise leads to larger angular occupancy for common colors**

391 This approximate theory is general in principle. It not only works for different environment priors,
392 but in principle should also work for the same environmental prior with different RNN's noise
393 strength (Figure 5). This scenario potentially mimics brains with different noise level.
394

395 We trained 50 RNNs under high noise and low noise conditions, respectively. Same procedures
396 were used to compute the RNNs attractors and angular occupancy. We found that high-noise RNNs
397 have larger angular occupancy to common colors (Figure 5A, B), and also more likely to form
398 attractors to represent common colors^{1,11} (Figure 5C, D).

399 To test whether the approximate formula (Equation 1) still works in this varying noise case, we
400 trained RNNs using different noise levels. Their dynamic dispersion, space occupancy and mean
401 bias correction were computed separately (Figure 5E, F, and G). Experimental memory error was
402 also computed by directly comparing input color to output color. Again, the approximate formula
403 aligns with experimental memory error well, and better than theory without considering the
404 decoding strategy (angular occupancy = 1) (Figure 5H). Note the error studied in Figure 5H is only
405 the error of the common colors. It can be observed that total error averaged across different input
406 colors naturally increases under larger noise conditions (Figure 5I).

407
408
409



410
411 **Figure 5. The results for different noise strength also support that angular occupancy is an important**
412 **factor to reduce memory error at common color.** (A) Decoding the delay plane for two example RNNs
413 trained on low-noise (noise strength $\sigma_{rec} = 0.1$) and high noise ($\sigma_{rec} = 0.3$). (B) Angular occupancy. Dash
414 vertical lines are common colors. Solid blue/red line indicates the mean of 50 RNNs and error bands are
415 standard deviation. (C) Neural dynamics during delay. Each trajectory is the delay of one trial with trajectory

416 color indicating the input color. Stars are the beginning of delay; dots are the ends. Black dots are attractors,
417 black crosses are saddles with long bar indicating the positive eigenvalue directions. **(D)** Number of
418 attractors for different colors. Attractors were concatenated from 50 RNNs. **(E, F, G)** Dynamic dispersion,
419 squared reciprocal of the angular occupancy and mean bias correction measured for 50 RNNs for each of
420 the noise strength (see Methods). Boxes indicate the interquartile range between the first and third quartiles
421 cross 50 RNNs with the central mark inside each box indicating the median. Whiskers extend to the lowest
422 and highest values within 1.5 times the interquartile range. **(H)** Comparing experimental results with
423 theoretical predictions for memory error. Dots are the median across 50 RNNs and error bar shows first and
424 third quantiles. In the experimental part, memory error of each RNN is the average of 5000 trials error fixing
425 input as common color (outlier trials' error were removed, see Methods). **(I)** Experimental error same as
426 **(H)**, but input colors were randomly sampled from an environmental prior $\sigma_s = 17.5^\circ$ instead of fixing to a
427 common color. In this figure, all RNNs were trained on environmental prior $\sigma_s = 17.5^\circ$. Noise was turned
428 off after training in computing panel **(A, B, C, D, F)**.

429 **Larger angular occupancy is due to both decoding dynamics and biased readout in the**
430 **decoding phase.**

431 We have shown that the decoding strategy (i.e., angular occupancy to a color) is an important
432 mechanism for reducing memory errors. This is an end-to-end effect of decoding phase. Next, we
433 reverse-engineered the decoding phase to explore what happens inside the decoding phase that led
434 to such a non-uniform angular occupancy. We study this problem by dividing and conquering. The
435 decoding phase can be roughly separated into three stages (Figure 3B, 6A). (1) Go Dynamics: a go
436 cue signal pushes the recurrent neural states from the delay plane to the response plane. (2)
437 Response Dynamics: recurrent neural states can dynamically move during the response epoch
438 (epoch length = 200 ms). (3) Readout: during the response epoch, recurrent neural states were
439 readout (Equation 3 in Methods) into response neural activity, then mapped into output colors
440 (Equation 7). The first two stages can be seen as continuing neural dynamics following the delay
441 epoch (although there is a go cue input during the go epoch, altering the dynamic trajectory). The
442 third stage, readout, does not involve dynamics. It simply multiplies the recurrent neural states with
443 a readout matrix, adds a bias current constant, and then transforms it into output color. Biologically
444 speaking, readout matrix may involve the neural connections from motor cortex to muscles (for
445 actions); bias current constant reflects the intrinsic property of each neuron (e.g. describing the
446 diversity of neural firing threshold). We inspect possible mechanisms leading to non-uniform
447 angular occupancy for each of these three stages (Go Dynamics, Response Dynamics, and Readout)
448 separately.

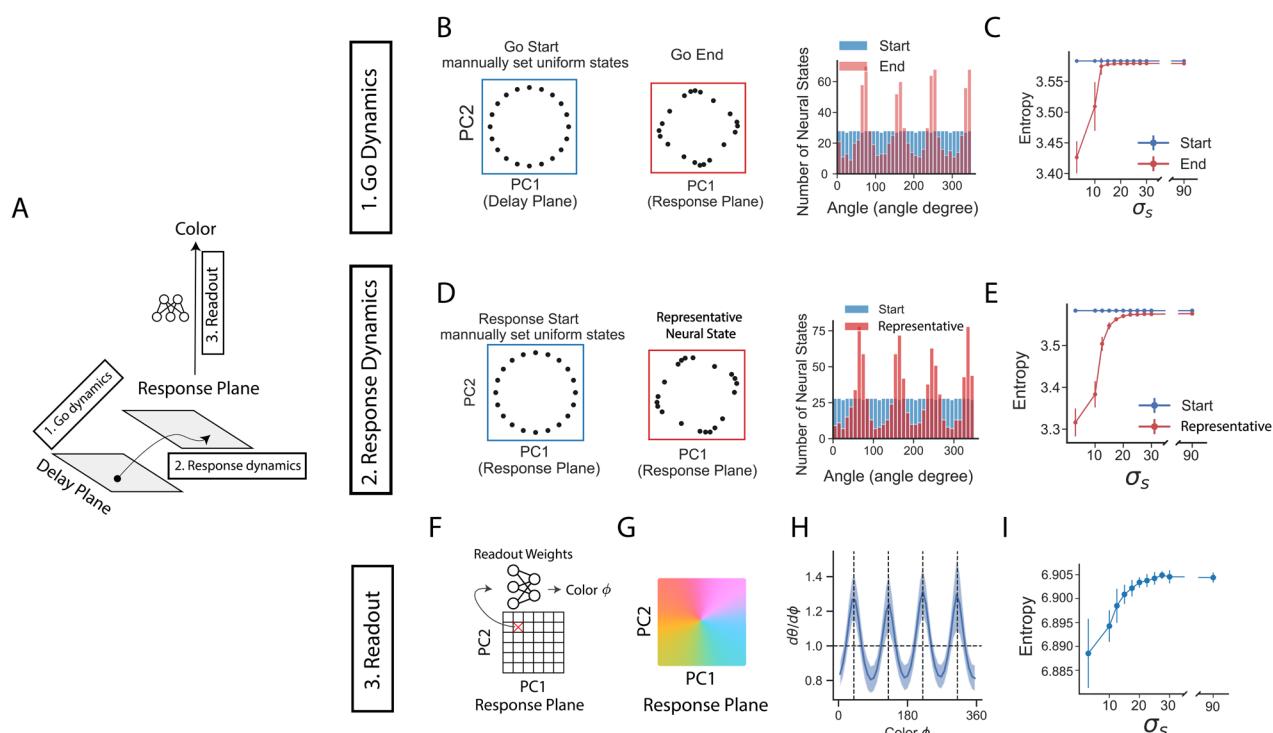


Figure 6. Larger angular occupancy (for common color) is due to both dynamics and biased readout mechanisms during the decoding phase. The single example RNN shown in this figure was trained on prior $\sigma_s = 3^\circ$; delay epoch length is fixed to 800 ms. **(A)** The decoding phase roughly contains three stages (see Figure 3A, B): (1) Go cue drives neural states from delay plane to response plane; (2) Neural dynamics during the response epoch; (3) Readout from recurrent neural states to output color. **(B)** We isolated the effect of Go Dynamics, by firstly fixing neural states uniformly along a ring in the delay plane, then letting the neural states evolve through go epoch. Left two boxes show one example RNN's neural states at the start and end of Go epoch. Right: the distribution of 1000 neural states of the same example RNN. **(C)** Entropy measures how non-uniform the neural states distribution is. Dots are the mean entropy across 50 RNNs (with outlier RNNs removed), error bar is the standard deviation. **(D, E)** Same as B, C, but the neural states were initialized in the response plane and evolved through response epoch. Representative neural states: recurrent neural states temporally averaged across whole response epoch, which is more relevant for decoding color (see Equation 7). **(F)** Decoding response plane, by selecting mesh points and then readout into color through a readout matrix and bias currents (see Methods, Equation 3 and 7). No dynamics were involved. **(G)** The example RNN's decoded response planes. **(H)** Angular occupancy in the response plane. Line shows the mean of 50 RNNs ($\sigma_s = 3^\circ$), and error band is the standard deviation. Vertical dash lines: four common colors. Outlier RNNs were removed. **(I)** Entropy of angular occupancy, dots are the mean and error bar are standard deviation (outlier RNNs removed).

In the Go Dynamics, tangent motion is a possible mechanism leading to larger angular occupancy for common color. We isolated the go epoch and measured its dynamics. At the beginning, we sampled a uniform ring of neural states (Figure 6B) in the delay plane, then let all these states evolve through go epoch (60 ms). At the end of the go epoch, we collected those neural states and fitted PC1-PC2 plane again. This refitted PC1-PC2 plane is called response plane. We found that neural states accumulated to a few positions. The degree of accumulation can be quantified by measuring the entropy of neural state distribution. The smaller the entropy is, the more likely the neural states tend to accumulate to a few positions. Figure 6C suggests that neural states accumulate happens when the environment prior σ_s is small. This biased motion during the go epoch is a possible reason leading to a non-uniform angular occupancy in the delay plane. For instance, if all neural states move to the same position during the go epoch, they will all be read out to the same color. This means that, regardless of where the delay neural states were, they will all be decoded into the same color. As a result, the angular occupancy is highly biased.

482

483 Similarly, biased motion during the response epoch may also contribute to non-uniform angular
484 occupancy. As in the study of Go Dynamics, we sampled uniform neural states on a ring in the
485 response plane and allowed them to evolve through the response epoch (200 ms). Since the
486 temporally averaged neural states were used to output a color, we computed the positions of
487 "representative neural states," which are the temporally averaged neural states during the response
488 epoch (Figure 6D). Results suggested that the tangent motion of neural states also occurs when the
489 environmental prior σ_s is small (Figure 6E). This implies that angular motion during the response
490 epoch could also be a reason for larger angular occupancy in the delay plane.
491

492

493 On the one hand, the above two dynamic mechanisms are reasonable. The go and response epochs
494 themselves require a duration of time, necessitating the maintenance of information throughout the
495 decoding phase, thus providing an opportunity for neural dynamics to play a role. However, on the
496 other hand, it is important to recognize that the dynamics during the go and response epochs can be
497 quite different from that of delay epoch. The go pulses acts as an input to the neural system,
498 prompting a shift in neural states during the go epoch. This shifted space (response plane) may have
499 different dynamic fields compared to the delay plane^{29,30}. Therefore, it should not be regarded as
500 merely a trivial extension of the delay dynamics.
501

502

503 Not only the dynamical evolution of recurrent state during the decoding phase matters, how the
504 recurrent neural state maps to an output color (including the mapping from recurrent neuron to
505 response neuron and population method for mapping response neural activities to an output color,
506 Equation 3 and 7) may be also important. This is called Readout and does not have a time factor
507 hence it is not a dynamics process. As an extreme example, a trivial Readout can map all possible
508 recurrent neural states (in the response plane) into a single color, hence the angular occupancy can
509 be highly biased.
510

511

512 To examine the Readout's characteristics, we sampled a grid mesh of points within the response
513 plane. Then, utilizing the readout matrix along with the bias current (as elaborated in Equation 3
514 and the population vector method in Equation 7), we transformed these mesh points into
515 corresponding colors (illustrated in Figure 6G). The resulting data (Figure 6H) revealed that the
516 Readout itself exhibits a bias, mapping larger angular spaces to common colors. This bias becomes
517 increasingly pronounced in a smaller environmental prior σ_s (Figure 6I).
518

519

520 Overall, these results suggest that angular motion during the go and response epoch, and the
521 biased readout process (including bias current in Equation 3) are the possible reasons for larger
522 angular occupancy for common color in the delay plane. Interestingly, we found that if RNNs
523 were trained on shorter response epochs, response dynamics became weaker while the readout
524 effect became stronger (SI Figure 3). This suggests that dynamic and non-dynamic factors may
525 compensate for each other. It also implies that in some scenarios, where the agent has limited time
526 to respond, the readout mechanism may play a more dominant role in shaping the decoding
527 strategy to reduce memory errors.
528

524

Discussion

525

526 Neural states and responses are noisy. How the noisy neural system maintains accurate information
527 is a central question in working memory. Previous researches emphasize neural mechanisms during
528 the delay. Our results provide a novel perspective, highlighting the importance of the decoding
529 phase in reducing memory errors. Smaller memory errors for a color can be achieved by allocating
530 larger state space occupancy (or angular occupancy in this paper). Further analysis revealed that
531 biased state space occupancy is due to both (1) continuing the attractor-based dynamics during the
532 decoding phase (although containing a go cue pulse) and (2) a biased readout from recurrent neural
533 activity to output color. These two factors mutually compensate as shown when training the RNNs
534 with shorter response epoch duration. In general, decoding phase (retrieval of information) is a
535 common phase of working memory tasks, our results encourage certain attention on the role of
536 decoding in diverse working memory tasks.
537

538

539 The logic of this paper resembles a previous human/animal behavioural experiment¹:
540 training/testing the agent to perform a memory task within a specific prior color distribution,
541 observing a reduction in memory errors for some input colors, and subsequently proposing and/or
542 testing potential neural mechanisms for this error reduction. But we conceive the suggested memory
543 error reduction mechanisms can be more universal than merely adapting to environmental priors.
544 For instance, in Figure 5, we demonstrated that RNNs trained on the same prior but with different
545 noise levels exhibit similar mechanisms for reducing memory errors. In a broader context,
546 mechanisms for reducing memory errors might arise when there is a high demand for accurate
547 retention of information. One potential example is reinforcement learning³¹. The importance of
548 information values is weighted by their associated rewards. Consequently, we further hypothesize
549 that RNNs may adopt similar mechanisms for reducing memory errors (attractor dynamics and a
550 larger state space occupancy) for high reward values. This hypothesis in future can possibly be
551 tested in physiology experiments^{1,30,32,33} by appropriately assigning rewards to different
552 information values while training animals to perform delayed-response tasks, or can possibly be
553 tested on RNN models using reinforcement learning^{34,35}.

554

555 Concretely, this paper proposes several experimental testable phenomena. First, when the input
556 value is important (i.e., high-prior color or high-reward value), the neural population dispersion
557 (measured by the variance of repeated trials) should be smaller. Second, state space occupancy for
558 important values should be larger. The state space occupancy can be measured similar to the
559 numerical procedure described here. Specifically, test the animal on multiple trials with different
560 input colors and collect the neural population states at the end of the delay. These neural states,
561 along with the animal's output color, provide a mapping from neural state space in the delay to the
562 output color. This mapping is determined by the decoding phase hence can be used for studying the
563 decoding strategy. With the mapping from neural state space to output color, numerical
564 differentiation can then be computed. State space occupancy for important colors (high-prior or
565 high-reward) can be compared with other colors. Finally, the approximate formula (Figure 4A-E)
566 can also be tested. For example, train the animal to perform delayed-response tasks in different
567 environmental priors^{1,36}, and then measure the dynamic dispersion, angular occupancy, and mean
568 bias correction separately. Theoretical predictions can be compared with the animal's experimental
569 memory error.

570

571 In certain physiology experiments, the neural population state resides within a low-dimensional
572 manifold and employs a ring-like structure to represent information^{13,30,37}, akin to the ring structure
573 discovered in our trained RNNs. In such cases, the state space occupancy can be computed by
determining the angle of the neural population state. However, we are aware that in many cases,

574 neural manifolds exist in high dimensions, such as a highly curved trajectory within a high-
575 dimensional space^{6,8}. In these instances, state space occupancy does not rely on angular occupancy
576 but rather possibly on “arc-length occupancy”. Arc-length occupancy means the arc length used to
577 represent a unit change of information. Measuring the arc-length occupancy of a highly curved
578 trajectory can be challenging. Yet this challenging can be eased by the concepts and methods
579 advanced in the recent geometric framework³⁸⁻⁴². Various nonlinear dimensional reduction
580 methods⁴² (e.g., CEBRA⁴³, UMAP⁴⁴) can project the manifold into a few latent dimensions for
581 better visualization. Additionally, methods like SPUD⁴⁵ or simple kernel smoothing³⁶ enable the
582 direct parameterization of a high-dimensional curve by its arc length, facilitating the measurement
583 of arc length occupancy for each information value. In future, it would be interesting to explore the
584 memory error reduction mechanisms in high-dimensional and more complex manifold scenarios.
585

586 We modeled a classical color-delayed response task^{9,11,13,16,17,46} that captures the most essential
587 three stages of memory process: perception, maintenance (delay), and decoding. Beyond this
588 classical experimental design, there are some other memory task variants^{1,2,30,36,47,48}. Some tasks
589 have no explicit go cues¹⁻³. For example, a color wheel can directly appear on the screen during the
590 response epoch to cue responses¹. This constant visual feedback to the animal may lead to neural
591 dynamics different from those shown in this paper. Furthermore, multiple items can be shown on
592 the screen for the animal to memorize⁴⁹, for example, simultaneously memorizing multiple colors³⁰.
593 How do the neural dynamics of multiple items interact with each other? Does the neural system
594 form a disentanglement representation for different items^{30,50,51}, or a shared mechanism for
595 correlated memory items⁵²? Is it possible that higher-rewarded items are encoded with larger
596 encoding spaces? Overall, the decoding phase (or so-called memory retrieval) exists ubiquitously
597 in all sorts of complex memory tasks. Our paper emphasizes the role of decoding in working
598 memory error reduction, and encourages future study of the decoding phase in more complex
599 memory tasks.
600

601 Beyond memory error reduction mechanisms, how RNNs adapt to the entire environmental prior is
602 an interesting problem. Adaptation is not simply about reducing the memory error for all colors. As
603 observed in Figure 1D, also in previous works^{1,9,10}, some low-prior colors exhibit larger memory
604 errors. Thus, there is a trade-off for the RNN: how much error should be reduced for high-prior
605 colors and how much should be increased for low-prior colors. This trade-off is unavoidable due to
606 the neural mechanism. Introducing attractors to reduce memory error for common colors
607 unavoidably leads to biased (or truncated) errors for nearby colors^{1,10}. In the angular occupancy
608 (Figure 3) part, enlarging common color’s angular occupancy will lead to smaller occupancy for
609 other colors. How do RNNs balance memory errors among different colors remains unclear.
610 Possible explanations can come from the Bayesian inference framework^{21,53}. For example, Sohn et
611 al.³⁶ assumes that the agent adapts to the environmental prior by making Bayesian inference in each
612 trial. Environmental prior is the prior function, noise distribution is the likelihood function. Based
613 on the prior and likelihood, they proposed that the agent needs to output a memorized information
614 value (in that paper was time duration) which equals to the mean of the posterior distribution. This
615 theory can make concrete prediction on the neural dynamics. Alternative theory suggests the
616 network can be a generative model for reproducing the environmental prior¹⁷. RNN trained in this
617 paper provides a ready-to-use models for theory pre-testing.

618 **Materials and Methods**

619 **RNN architecture**

620 The RNN consists of 12 perception input neurons, 1 go input neuron, 12 response output neurons
621 and 256 fully connected recurrent neurons (Fig. 1A). The dynamic equation of the RNN is⁵⁴

$$\mathbf{x}_t = (1 - \alpha)\mathbf{x}_{t-1} + \alpha \left(\mathbf{W}^{rec} \mathbf{r}_{t-1}^{rec} + \mathbf{W}^{in} \mathbf{u}_t + \sqrt{2\alpha^{-1}\sigma_{rec}^2} \boldsymbol{\epsilon} + \mathbf{b} \right), \quad (2)$$

622 where \mathbf{x}_t is a 256-dimensional vector representing the activities/state of recurrent neurons (neural
623 state) at time step t , and $\alpha = \Delta t / \tau$, where Δt is the time length for each computational step, τ is
624 the neuron membrane time constant. In this study, we set $\Delta t = 20$ ms and $\alpha = 1$, which is the same
625 choice as in previous work^{25,54,55}. The change of neural state \mathbf{x}_{t+1} subjects to 4 factors: (1) The
626 inputs from the other recurrent neurons, $\mathbf{W}^{rec} \mathbf{r}_{t-1}^{rec}$, where $\mathbf{r}_t^{rec} = \tanh(\mathbf{x}_t)$ and \mathbf{W}^{rec} is the
627 recurrent weights but the diagonal elements fixed to 0 (self-connections are not allowed)⁵⁴ (2) The
628 inputs from perception and go neurons, $\mathbf{W}^{in} \mathbf{u}_t$, where \mathbf{u}_t consists of 12 components for the
629 perception neurons and 1 component for the go neuron, and \mathbf{W}^{in} is the input weights. During the
630 perception epoch, the perception neurons is entangled with an addictive gaussian noise with
631 standard deviation $\sigma_x = 0.2$. (3) The intrinsic recurrent noise in the RNN, $\sqrt{2\alpha^{-1}\sigma_{rec}^2} \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ is
632 a 256-dimensional noisy vector drawn from the standard normal distribution at each time step, and
633 σ_{rec} is the recurrent neural intrinsic noise strength. $\sigma_{rec} = 0.2$ was used in all main figures, except
634 in Figure 5 where we inspected different noise strengths. (4) A constant bias current vector \mathbf{b} .

635
636 The activities of the 12 response neurons are read from the \mathbf{r}_t^{rec} ,

$$\mathbf{z}_t = \mathbf{W}^{out} \mathbf{r}_t^{rec} + \mathbf{b}^{out}. \quad (3)$$

637 **A delayed-response trial and RNN training**

638 A trial contained a fixation, perception, delay, go, and response epochs. During the fixation epoch,
639 the RNN didn't receive any inputs (i.e., $\mathbf{u}_t = 0$). In the perception epoch, a color ϕ_s was drawn
640 from a prior distribution,

$$P_s(\phi_s) = \frac{1}{4} \sum_i \mathcal{VM}(\phi_s - \mu_i; \sigma_s^2), \quad (4)$$

641 where $VM(\cdot)$ is a von Mises function (circular version of the normal distribution) with mean values
642 at $\mu_i \in \{40^\circ, 130^\circ, 220^\circ, 310^\circ\}$

$$\mathcal{VM}(\phi_s - \mu_i; \sigma_s^2) = \frac{1}{2\pi I_0(1/\sigma_s^2)} \exp[\cos(\phi_s - \mu_i)/\sigma_s^2], \quad (5)$$

643 where σ_s is a free parameter to control the “width” of the prior distribution, and $I_0(\cdot)$ is the modified
644 Bessel function of order 0 for normalization purpose. Units of ϕ_s , μ and σ_s were converted from
645 degree to rad before feeding into the von Mises function. A color was drawn from the above prior
646 distribution and then was sensed by the 12 perception neurons. Perception neurons had tuning
647 curves that were also von Mises functions $\mathcal{VM}(\phi_s - \mu_{p,i}; \sigma_p^2)$ with $\mu_{p,i} = \frac{360^\circ}{12} i$ and $\sigma_p = \frac{1}{2} \frac{360^\circ}{12}$
648 for $i = 0, 1, \dots, 11$ (Fig. 1A). Therefore, in the perception epoch, the perception components of the

649 input vector \mathbf{u}_t are $u_{t,i} = \mathcal{VM}(\theta_s - \mu_{p,i}; \sigma_p^2) + \epsilon_x$ where ϵ_x is a gaussian noise with standard
650 deviation $\sigma_x = 0.2$, and the go component is 0 ($u_{t,12} = 0$).
651

652 Next, during the delay epoch, all inputs were 0. The delay length was drawn from an uniform
653 distribution $U(0 \text{ ms}, 1000 \text{ ms})$ in each trial. Then, in the go epoch, we set $u_{t,12} = 1$ to indicate a go
654 cue and the rest of inputs and the outputs were 0. Finally, in the response epoch, the 12 response
655 neurons were required to reproduce the neural activity of same as the perception neurons during
656 their perception epoch, i.e. $\tilde{\mathbf{z}}_{t,i} = \mathcal{VM}(\phi_s - \mu_{p,i}; \sigma_p^2)$ for $i = 0, 1, \dots, 11$.
657

658 The loss function is

$$\text{Loss} = \frac{1}{T_{tot}} \sum_t m_t \left[|\mathbf{z}_t - \tilde{\mathbf{z}}_t|^2 + \frac{1}{N_{rec}} (\beta T |W_{rec}|^2 + \gamma |r_{rec}(x) + 1|^2) \right], \quad (6)$$

659 where m_t is the mask that is 0 in the fixation epoch and 1 in the rest epochs, T_{tot} is the trial time
660 length, N_{rec} is the number of recurrent neurons, β and γ in the regularization terms constrain the
661 recurrent weights and neural firing rates. $\tilde{\mathbf{z}}_t$ is the target response output. It was zero across
662 perception, delay and go epochs, and was the target response $\tilde{\mathbf{z}}_{t,i} = \mathcal{VM}(\phi_s - \mu_{p,i}; \sigma_p^2)$ during
663 the response epoch. This loss function was minimized by gradient descent using the Adams
664 optimizer⁵⁶.
665

666 Directly training the RNN to optimize the above loss function is difficult, we employed the
667 progressive training protocol²⁷: First, we trained the RNN with no noise, no regularization terms,
668 uniform prior distribution, and delay epoch length fixed to 0. Secondly, we retrained the RNN with
669 delay epoch length drawn from $U(0 \text{ ms}, 1000 \text{ ms})$. Thirdly, the RNN is retrained considering the
670 noise and regularizations. The resulting model is called the pretrained RNN. Finally, we retrained
671 the pretrained RNN with the desired prior distribution. If the desired prior distribution is still a
672 uniform prior, then the trained RNN was called as Uniform RNN (approximately equivalent to the
673 prior $\sigma_s = 90^\circ$), otherwise was called as Biased RNN.

674 Mapping response neural activity into color

675 The actual response neural activity can be mapped back to a color by a population method. The
676 mapped color is also called the RNN's output color. More precisely, the response neural activity \mathbf{z}_t
677 was firstly taking averaged within middle time window of the response epoch (60 ms to 140 ms,
678 note the total response epoch is 200 ms). Next, the averaged response neural activity $\bar{\mathbf{z}}$ was
679 converted into the output color by a populational vector method^{54,57}

$$\hat{\phi} = \text{ang} \left(\sum_{m=1}^{N=12} \bar{z}_m e^{i\mu_{p,m}} \right), \quad (7)$$

680 where $\text{ang}(\cdot)$ means taking the angle of a complex number, i here is a unit imaginary number, and
681 \bar{z}_m is the m th components of vector $\bar{\mathbf{z}}$, i.e. the m th response neuron. $\mu_{p,m} = \frac{360^\circ}{12} m$ which was set
682 to be align with the perception tuning mean.

683 Removing outlier RNNs

684 Due to the stochastic nature of the training process, some RNNs exhibited significantly different
685 metric values (e.g., memory error, entropy etc.) compared to others, this may be due to a failure of
686 training. To mitigate the effect of these outliers, this paper employs a criterion for exclusion: any
687 RNN whose measured metric exceeds 1.5 times the interquartile range (IQR) above the third
688 quartile or below the first quartile is removed from the analysis. This exclusion was specified by
689 the phrase 'outlier RNNs were removed' in figure legends.

690 **Identifying fix points**

691 The dynamic equation of the RNN, Eq. (1), can be reformulated in the form of

$$\mathbf{x}_t - \mathbf{x}_{t-1} = F(\mathbf{x}_{t-1}), \quad (8)$$

692 which defines the velocity field $\mathbf{v}(\mathbf{x}) = F(\mathbf{x}) / \Delta t$. In this study, we set $\Delta t = 20$ ms. The global
693 minimum points of speed $v(\mathbf{x}) = |\mathbf{v}(\mathbf{x})| = 0$ are the fixed points. Specifically, we ran the trained
694 RNN with 500 trials with stimuli equally distributed from 0° to 360° . The neural states at the
695 beginning of delay epochs were collected. According to the visualization in Figure 3, these neural
696 states were likely to near the ring on the delay. Hence these neural states were used as the
697 initialization of searching point. Starting from these initial neural states, we searched for the global
698 minima of speed using the gradient descent, where the independent variable is the position \mathbf{x} . Due
699 to the imperfection of the gradient descent algorithm, we may end up with local minima.
700 Nevertheless, these local minimum points are similar to the fixed points⁵⁸, hence they are also called
701 fixed points in this paper. For each fixed point, we computed the eigenvalues and eigenvectors of
702 its Jacobian matrix $J = \frac{\partial F(\mathbf{x})}{\partial \mathbf{x}}$. A fixed point is called an attractor if its largest eigenvalue is negative,
703 while it is called saddle point if some eigenvalues are positive. The eigenvectors of the positive
704 eigenvalues are projected and shown on the PC1-PC2 plane (Fig. 3). Only fixed points near the
705 delay ring were shown.

706 **Decoding an end-of-delay neural state**

707 We built an RNN decoder to decode a delay neural state \mathbf{x}_t^{rec} to its corresponding color. This RNN
708 decoder faithfully represents the decoding phase of the RNN. Specifically, an RNN decoder is the
709 same as the original trained RNN, but its recurrent neural state is initialized as the to-be-decoded
710 neural state. Then, the RNN decoder runs directly through the go and response epoch. Response
711 neural activities in the response epoch were temporally averaged. Averaged response was then
712 decoded to an output color by the population vector method (Equation 7). This RNN decoder was
713 used, for example, to decode attractors color (Figure 3D).

714
715 We also used the RNN decoder to decode the PC1-PC2 plane of the delay. A RNN was firstly ran
716 through 1000 trials with uniformly distributed input color. Neural states during the delay were
717 collected, used for fitting the PC1-PC2 plane, also called the delay plane. Next, a grid of mesh
718 points (50 times 50) was sampled in the delay plane, centered at the delay plane center. Each of
719 these mesh points were then PCA inversed back the original high-dimensional space, finally
720 decoded by the RNN decoder (Figure 3G).

721
722 We also used the RNN decoder to find the relation between the angle of the delay plane and the
723 color. Similar as above, a RNN was firstly ran through 1000 trials with uniformly distributed input
724 color. Neural states at the end of delay (delay length is fixed as 800 ms) were collected, for fitting
725 the PC1-PC2 plane (delay plane). Next, within the delay plane, a radius was computed, which was
726 defined as the averaging distances of all collected 1000 neural states to the delay plane center. Using

727 this radius, dense points (1000 points) of a ring were resampled. Finally, each of the resample points
728 were decoded by the RNN decoder. This provides a mapping from the angle of the delay plane to
729 the color (Figure 3H).

730 **Cross-decoding experiment**

731 Cross-decoding involves preparing delay neural states using one RNN model (model 1) and
732 decoding the prepared delay neural state using another RNN model (model 2). Specifically, model
733 1 was randomly chosen from either the Biased RNN or Uniform RNN categories. This model
734 underwent 500 trials with a fixed input of a common color. Neural states at the end of delay were
735 collected. Next, these end-of-delay neural states were used to set the recurrent neural states of model
736 2 which was randomly selected from the second category (which may be the same as model 1's
737 category). Model 2 processed these states over its decoding phase, outputted colors. The squared-
738 root of the mean squared error across 500 trials was defined as the memory error. This procedure
739 was repeated to gather memory error data for 50 pairs from each category combination, namely
740 Bias & Bias, Bias & Uniform, Uniform & Bias, and Uniform & Uniform. The first category
741 represents model 1 while the second represents model 2.
742

743 A key challenge in this process is the “neuron-matching problem” which involves determining how
744 to match neurons in model 2 with those in model 1. Our approach to neuron matching is based on
745 comparing the neural preferred color rank. The preferred color for a neuron in an RNN model is
746 determined as the following. Dense neural states along a ring on the PC1-PC2 delay plane were
747 sampled; then were decoded into output colors. This provides a delay-neural-states-to-color
748 mapping. Since delay neural state is nothing but a vector of concatenated individual neural
749 activations, delay-neural-states-to-color mapping was also serves as a neural-activation-to-color
750 mapping. The preferred color for each neuron was defined as the color that maximized its neural
751 activation.
752

753 We computed the preferred colors of all recurrent neurons in Model 1 and ranked them from 0
754 degrees (smallest) to 360 degrees (largest). We repeated this process for Model 2. We achieved
755 neuron pairing by matching these ranks. Therefore, the prepared delay neuron state in Model 1,
756 which includes 256 neurons, can correspond one-to-one with the recurrent neural state in Model 2,
757 which also comprises 256 neurons. Subsequently, model 2 underwent the decoding phase to
758 compute the memory error.

759 **A Taylor expansion approximation to the memory error**

760 The squared memory error of a fixed common color ϕ_c is

$$\epsilon^2(\phi_c) = \sum_{i=1}^N (\phi_i - \phi_c)^2 / N, \quad (9)$$

761 where N is the number of trials, ϕ_i is the RNN's reported color in trial i . Approximately the RNN
762 encodes a color by the angle (θ) of the delay plane, so output color is a function of the neural state's
763 angle $\phi(\theta)$. Denoting the angle encoding common color as θ_c . Assuming during delay the actual
764 neural state's angle doesn't deviate from θ_c too much, then we can Taylor expand $\phi(\theta) \approx \phi_c +$
765 $\frac{d\phi}{d\theta}|_{\phi_c}(\theta - \theta_c)$. Substituting this back to Equation (9) yields
766

$$\begin{aligned}\epsilon^2(\phi_c) &= \sum_{i=1}^N (\phi_i - \phi_c)^2 / N \\ &\approx \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 \sum_{i=1}^N (\theta_i - \theta_c)^2 / N \\ &= \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 \sum_{i=1}^N [(\theta_i - \bar{\theta}) + (\bar{\theta} - \theta_c)]^2 / N \\ &= \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 \sum_{i=1}^N \frac{(\theta_i - \bar{\theta})^2}{N} + \left(\frac{d\phi}{d\theta} \Big|_{\phi_c} \right)^2 (\bar{\theta} - \theta_c)^2,\end{aligned}\tag{10}$$

were $d\phi/d\theta|_{\phi_c}$ is the reciprocal of the angular occupancy at common color, $\sum_{i=1}^N (\theta_i - \bar{\theta})^2 / N$ is the dynamic dispersion, and the last term describes a correction from the actually neural state mean angle to the common color angle.

770 Numerical procedure for testing the Taylor expansion approximation.

To test this approximation formula (Equation 10), we computed the dynamic dispersion, angular occupancy and mean bias correction, respectively. To compute the dynamic dispersion, we ran each of the RNN 5000 trials with fixed input color as common color (delay length 800 ms). Neural states at the end of the delay were collected, and their corresponding angles θ_i were computed. Outlier angles (1.5 IQR below the first quantile or above the third quantile) were removed. Dynamic dispersion was measured as the variance of θ_i . The mean of θ_i was also denoted as $\bar{\theta}$. To compute angular occupancy, similar as Figure 3H (see Methods), dense points along a ring were sampled in the delay plane. Each of these points was decoded by the RNN decoder, hence we obtained the angle-color mapping. Numerical differentiation was then performed to compute the angular occupancy. In addition, the angle on the ring whose color was closest to the common color was denoted as θ_c . Finally, using the $\bar{\theta}$ and θ_c computed previously, mean bias correction can be computed as written in the Equation (10).

The resulting memory error computed from all these terms (equation 10) was called theoretical prediction. To explore the importance of decoding strategy (i.e. angular occupancy), as a comparison, we also computed theoretical prediction but setting angular occupancy to 1.

For the experimental RNN's memory error, for each condition (e.g. environmental prior σ_s or noise level σ_{rec}), 50 RNNs were used. Each RNN was ran on 5000 trials. After removing the outliers' trial errors, each RNN's experimental memory error is the trial error's standard deviation (squared root of Equation 9).

792 Decoding neural states on the response plane

A neural state on the delay plane can be decoded into an output color by using the RNN decoder which is continuing running the RNN through go and response epochs. Similarly, a neural state in the response plane can also be decoded into an output color by continuing running the RNN. Since the response is already the last epoch, hence the "continuously running" only means reading the (recurrent) neural state out to the response neuron activity, by using RNN's readout weight \mathbf{W}^{out}

798 and bias current \mathbf{b}^{out} (Equation 3). The response neuron activity was then translated into an output
799 color using populational vector method (equation 7). Note that no dynamics is involved since no
800 time component here. This decoding procedure is also called readout decoder. Readout decoder's
801 properties is fully determined by the readout weight and bias currents.
802

803 To inspect the properties of Readout process, we used the Readout decoder to decode mesh points
804 on the response PC1-PC2 plane. Specifically, each RNN was run through 1000 trials, neural states
805 during the response epochs were collected for fitting the response plane (PC1-PC2 plane). After
806 fitting, we mesh grid points on the response plane. They were further decoded by the Readout
807 decoder.
808

809 Similarly, we can also inspect the relation between angle and the color in the response plane. The
810 procedure is analogous to what was described in the “Decoding an end-of-delay neural state”
811 session, except that the delay plane was replaced by the response plane, and readout decoders were
812 used instead of RNN decoders.

813 **References**

- 814 1. Panichello, M. F., DePasquale, B., Pillow, J. W. & Buschman, T. J. Error-correcting
815 dynamics in visual working memory. *Nat. Commun.* **10**, 1–11 (2019).
- 816 2. Lundqvist, M., Rose, J., Brincat, S. L., Warden, M. R., Buschman, T. J., Herman, P. &
817 Miller, E. K. Reduced variability of bursting activity during working memory. *Sci. Rep.* **12**,
818 1–10 (2022).
- 819 3. Yu, Q., Panichello, M. F., Cai, Y., Postle, B. R. & Buschman, T. J. Delay-period activity in
820 frontal, parietal, and occipital cortex tracks noise and biases in visual working memory.
821 *PLoS Biol.* **18**, 1–17 (2020).
- 822 4. Melton, A. W. Implications of short-term memory for a general theory of memory. *J.
823 Verbal Learning Verbal Behav.* **2**, 1–21 (1963).
- 824 5. Burak, Y. & Fiete, I. R. Fundamental limits on persistent activity in networks of noisy
825 neurons. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 17645–17650 (2012).
- 826 6. Khona, M. & Fiete, I. R. Attractor and integrator networks in the brain. *Nat. Rev. Neurosci.*
827 **23**, 744–766 (2022).
- 828 7. Renart, A., Song, P. & Wang, X. J. Robust spatial working memory through homeostatic
829 synaptic scaling in heterogeneous cortical networks. *Neuron* **38**, 473–485 (2003).
- 830 8. Shaham, N. & Burak, Y. Slow diffusive dynamics in a chaotic balanced neural network.
831 *PLoS Comput. Biol.* **13**, 1–26 (2017).
- 832 9. Eissa, T. L. & Kilpatrick, Z. P. Learning efficient representations of environmental priors
833 in working memory. *PLoS Comput. Biol.* **19**, (2023).
- 834 10. Kilpatrick, Z. P., Ermentrout, B. & Doiron, B. Optimizing working memory with
835 heterogeneity of recurrent cortical excitation. *J. Neurosci.* **33**, 18999–19011 (2013).
- 836 11. Xiong, H.-D. & Wei, X.-X. Optimal encoding of prior information in noisy working
837 memory systems. *Conf. Cogn. Comput. Neurosci. (CCN, 2022)* (2022)
doi:10.32470/ccn.2022.1162-0.
- 838 12. Chaudhuri, R. & Fiete, I. Computational principles of memory. *Nat. Neurosci.* **19**, 394–403
839 (2016).
- 840 13. Murray, J. D., Bernacchia, A., Roy, N. A., Constantinidis, C., Romo, R. & Wang, X. J.
841 Stable population coding for working memory coexists with heterogeneous neural
842 dynamics in prefrontal cortex. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 394–399 (2017).
- 843 14. Seeholzer, A., Deger, M. & Gerstner, W. *Stability of working memory in continuous
844 attractor networks under the control of shortterm plasticity*. *PLoS Computational Biology*
845 vol. 15 (2019).
- 846 15. Inagaki, H. K., Fontolan, L., Romani, S. & Svoboda, K. Discrete attractor dynamics
847 underlies persistent activity in the frontal cortex. *Nature* **566**, 212+ (2019).
- 848 16. Darshan, R. & Rivkind, A. Learning to represent continuous variables in heterogeneous
849 neural networks. *Cell Rep.* **39**, 110612 (2022).
- 850 17. Kilpatrick, Z. P. Synaptic mechanisms of interference in working memory. *Sci. Rep.* **8**, 1–
851 20 (2018).
- 852 18. Saxena, S. & Cunningham, J. P. Towards the neural population doctrine. *Curr. Opin.
853 Neurobiol.* **55**, 103–111 (2019).
- 854 19. Haefner, R. M., Gerwinn, S., MacKe, J. H. & Bethge, M. Inferring decoding strategies
855 from choice probabilities in the presence of correlated variability. *Nat. Neurosci.* **16**, 235–
856 242 (2013).
- 857 20. McGinty, V. B. & Lupkin, S. M. Behavioral read-out from population value signals in
858 primate orbitofrontal cortex. *Nat. Neurosci.* **26**, (2023).
- 859 21. Fritzsche, M., Spaak, E. & de Lange, F. P. A bayesian and efficient observer model explains
860 concurrent attractive and repulsive history biases in visual perception. *Elife* **9**, 1–32 (2020).
- 861 22. Ni, A. M., Huang, C., Doiron, B. & Cohen, M. R. A general decoding strategy explains the

- 863 relationship between behavior and correlated variability. *Elife* **11**, 1–21 (2022).
- 864 23. Sadtler, P. T., Quick, K. M., Golub, M. D., Chase, S. M., Ryu, S. I., Tyler-Kabara, E. C.,
865 Yu, B. M. & Batista, A. P. Neural constraints on learning. *Nature* **512**, 423–426 (2014).
- 866 24. Williams, A. H., Kim, T. H., Wang, F., Vyas, S., Ryu, S. I., Shenoy, K. V., Schnitzer, M.,
867 Kolda, T. G. & Ganguli, S. Unsupervised Discovery of Demixed, Low-Dimensional Neural
868 Dynamics across Multiple Timescales through Tensor Component Analysis. *Neuron* **98**,
869 1099–1115.e8 (2018).
- 870 25. Orhan, A. E. & Ma, W. J. A diverse range of factors affect the nature of neural
871 representations underlying short-term memory. *Nat. Neurosci.* **22**, 275–283 (2019).
- 872 26. Yang, G. R. & Wang, X. J. Artificial Neural Networks for Neuroscientists: A Primer.
873 *Neuron* **107**, 1048–1070 (2020).
- 874 27. Chaisangmongkon, W., Swaminathan, S. K., Freedman, D. J. & Wang, X. J. Computing by
875 Robust Transience: How the Fronto-Parietal Network Performs Sequential, Category-
876 Based Decisions. *Neuron* **93**, 1504–1517.e4 (2017).
- 877 28. Bae, G., Olkkonen, M., Allred, S. R. & Flombaum, J. I. Why Some Colors Appear More
878 Memorable Than Others : A Model Combining Categories and Particulars in Color
879 Working Memory. *J. Exp. Psychol. Gen.* **144**, 744–763 (2015).
- 880 29. Smith, J. T. H., Linderman, S. W. & Sussillo, D. Reverse engineering recurrent neural
881 networks with Jacobian switching linear dynamical systems. *Adv. Neural Inf. Process. Syst.*
882 **20**, 16700–16713 (2021).
- 883 30. Panichello, M. F. & Buschman, T. J. Shared mechanisms underlie the control of working
884 memory and attention. *Nature* **592**, 601–605 (2021).
- 885 31. Botvinick, M., Wang, J. X., Dabney, W., Miller, K. J. & Kurth-Nelson, Z. Deep
886 Reinforcement Learning and Its Neuroscientific Implications. *Neuron* vol. 107 603–616
887 (2020).
- 888 32. Constantinidis, C., Franowicz, M. N. & Goldman-Rakic, P. S. Coding specificity in cortical
889 microcircuits: A multiple-electrode analysis of primate prefrontal cortex. *J. Neurosci.* **21**,
890 3646–3655 (2001).
- 891 33. Funahashi, S., Bruce, C. J. & Goldman-Rakic, P. S. Mnemonic coding of visual space in
892 the monkey's dorsolateral prefrontal cortex. *J. Neurophysiol.* **61**, 331–349 (1989).
- 893 34. Song, H. F., Yang, G. R. & Wang, X. J. Reward-based training of recurrent neural
894 networks for cognitive and value-based tasks. *Elife* **6**, 1–24 (2017).
- 895 35. Wang, J. X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z.,
896 Hassabis, D. & Botvinick, M. Prefrontal cortex as a meta-reinforcement learning system.
897 *Nat. Neurosci.* **21**, 860–868 (2018).
- 898 36. Sohn, H., Narain, D., Meirhaeghe, N. & Jazayeri, M. Bayesian Computation through
899 Cortical Latent Dynamics. *Neuron* **103**, 934–947.e5 (2019).
- 900 37. Wolff, M. J., Jochim, J., Akyürek, E. G., Buschman, T. J. & Stokes, M. G. Drifting codes
901 within a stable coding scheme for working memory. *PLoS Biol.* **18**, 1–19 (2020).
- 902 38. DiCarlo, J. J. & Cox, D. D. Untangling invariant object recognition. *Trends Cogn. Sci.* **11**,
903 333–341 (2007).
- 904 39. Chung, S. Y. & Abbott, L. F. Neural population geometry: An approach for understanding
905 biological and artificial neural networks. *Curr. Opin. Neurobiol.* **70**, 137–144 (2021).
- 906 40. Jazayeri, M. & Ostojevic, S. Interpreting neural computations by examining intrinsic and
907 embedding dimensionality of neural activity. *Curr. Opin. Neurobiol.* **70**, 113–120 (2021).
- 908 41. Barack, D. L. & Krakauer, J. W. Two views on the cognitive brain. *Nat. Rev. Neurosci.* **22**,
909 359–371 (2021).
- 910 42. Cunningham, J. P. & Yu, B. M. Dimensionality reduction for large-scale neural recordings.
911 *Nat. Neurosci.* **17**, 1500–1509 (2014).
- 912 43. Schneider, S., Lee, J. H. & Mathis, M. W. Learnable latent embeddings for joint

- 913 behavioural and neural analysis. *Nature* **617**, 360–368 (2023).
- 914 44. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and
915 Projection for Dimension Reduction. *arXiv* (2018).
- 916 45. Chaudhuri, R., Gercek, B., Pandey, B., Peyrache, A. & Fiete, I. The intrinsic attractor
917 manifold and population dynamics of a canonical cognitive circuit across waking and
918 sleep. *Nat. Neurosci.* **22**, 1512–1520 (2019).
- 919 46. Roux, F. & Uhlhaas, P. J. Working memory and neural oscillations: Alpha-gamma versus
920 theta-gamma codes for distinct WM information? *Trends Cogn. Sci.* **18**, 16–25 (2014).
- 921 47. Apostel, A., Panichello, M., Buschman, T. J. & Rose, J. Corvids optimize working memory
922 by categorizing continuous stimuli. *Commun. Biol.* **6**, 1–13 (2023).
- 923 48. Mou, X., Suresh, P. & Ji, D. Tetrode recording of rat CA1 place cells in an observational
924 spatial working memory task. *STAR Protoc.* **3**, 101501 (2022).
- 925 49. Lara, A. H. & Wallis, J. D. Executive control processes underlying multi-item working
926 memory. *Nat. Neurosci.* **17**, 876–883 (2014).
- 927 50. Boyle, L., Posani, L., Irfan, S., Siegelbaum, S. A. & Fusi, S. The geometry of hippocampal
928 CA2 representations enables abstract coding of social familiarity and identity. *bioRxiv*
929 2022.01.24.477361 (2022).
- 930 51. Flesch, T., Juechems, K., Dumbalska, T., Saxe, A. & Summerfield, C. Orthogonal
931 representations for robust context-dependent task performance in brains and neural
932 networks. *Neuron* 1–13 (2022) doi:10.1016/j.neuron.2022.01.005.
- 933 52. Al Roumi, F., Marti, S., Wang, L., Amalric, M. & Dehaene, S. Mental compression of
934 spatial sequences in human working memory using numerical and geometrical primitives.
935 *Neuron* **109**, 2627-2639.e4 (2021).
- 936 53. Lange, R. D., Shivkumar, S., Chattoraj, A. & Haefner, R. M. Bayesian encoding and
937 decoding as distinct perspectives on neural coding. *Nat. Neurosci.* **26**, 2063–2072 (2023).
- 938 54. Bi, Z. & Zhou, C. Understanding the computation of time using neural network models.
939 *Proc. Natl. Acad. Sci. U. S. A.* **117**, 10530–10540 (2020).
- 940 55. Wang, J., Narain, D., Hosseini, E. A. & Jazayeri, M. Flexible timing by temporal scaling of
941 cortical responses. *Nat. Neurosci.* **21**, 102–112 (2018).
- 942 56. Kingma, D. P. & Ba, J. L. Adam: A method for stochastic optimization. *3rd Int. Conf.
943 Learn. Represent. ICLR 2015 - Conf. Track Proc.* 1–15 (2015).
- 944 57. Wimmer, K., Nykamp, D. Q., Constantinidis, C. & Compte, A. Bump attractor dynamics in
945 prefrontal cortex explains behavioral precision in spatial working memory. *Nat. Neurosci.*
946 **17**, 431–439 (2014).
- 947 58. Sussillo, D. & Barak, O. Opening the black box: Low-dimensional dynamics in high-
948 dimensional recurrent neural networks. *Neural Comput.* **25**, 626–649 (2013).
- 949
- 950
- 951

952 Acknowledgments:

953 We thank Daoyun Ji from the Baylor College of Medicine for constructive comments and
954 suggestions to this work.

956 Funding:

957 Hong Kong Baptist University Strategic Development Fund
958 the Hong Kong Research Grant Council (GRF12200620)
959 the Hong Kong Baptist University Research Committee Interdisciplinary Research
960 Clusters Matching Scheme 2018/19 (RC-IRCMs/18 19/SCI01)
961 the National Science Foundation of China (Grant 11975194)

963 **Author contributions:**

964 Conceptualization: ZY, LT, CZ
965 Methodology: ZY, HL, LT, CZ
966 Investigation: ZY, HL
967 Supervision: LT, CZ
968 Writing—original draft: ZY, HL, LT, CZ
969 Writing—review & editing: ZY, HL, LT, CZ
970

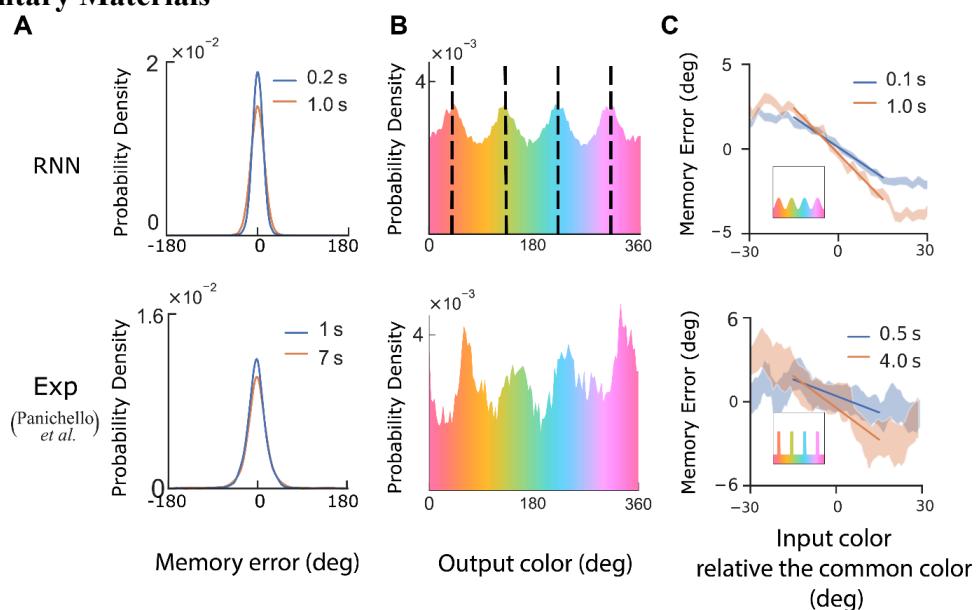
971 **Competing interests:** Authors declare that they have no competing interests
972

973 **Declaration of generative AI and AI-assisted technologies in the writing process:** During the
974 preparation of this work the author(s) used ChatGPT-4 in order to polish the language and
975 improve the readability of the manuscript. After using this tool/service, the authors reviewed and
976 edited the content as needed and take full responsibility for the content of the published article.
977

978 **Data and materials availability:** The codes for training and analysis are available at
979 https://github.com/AgeYY/working_memory.git
980

981
982

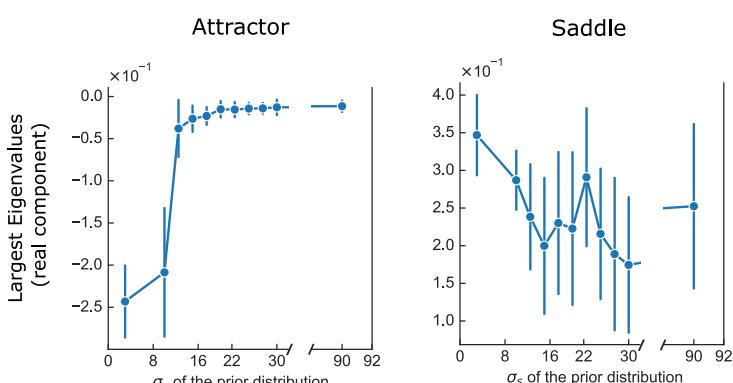
Supplementary Materials



983
984
985
986
987
988
989
990
991
992
993
994
995
996
997

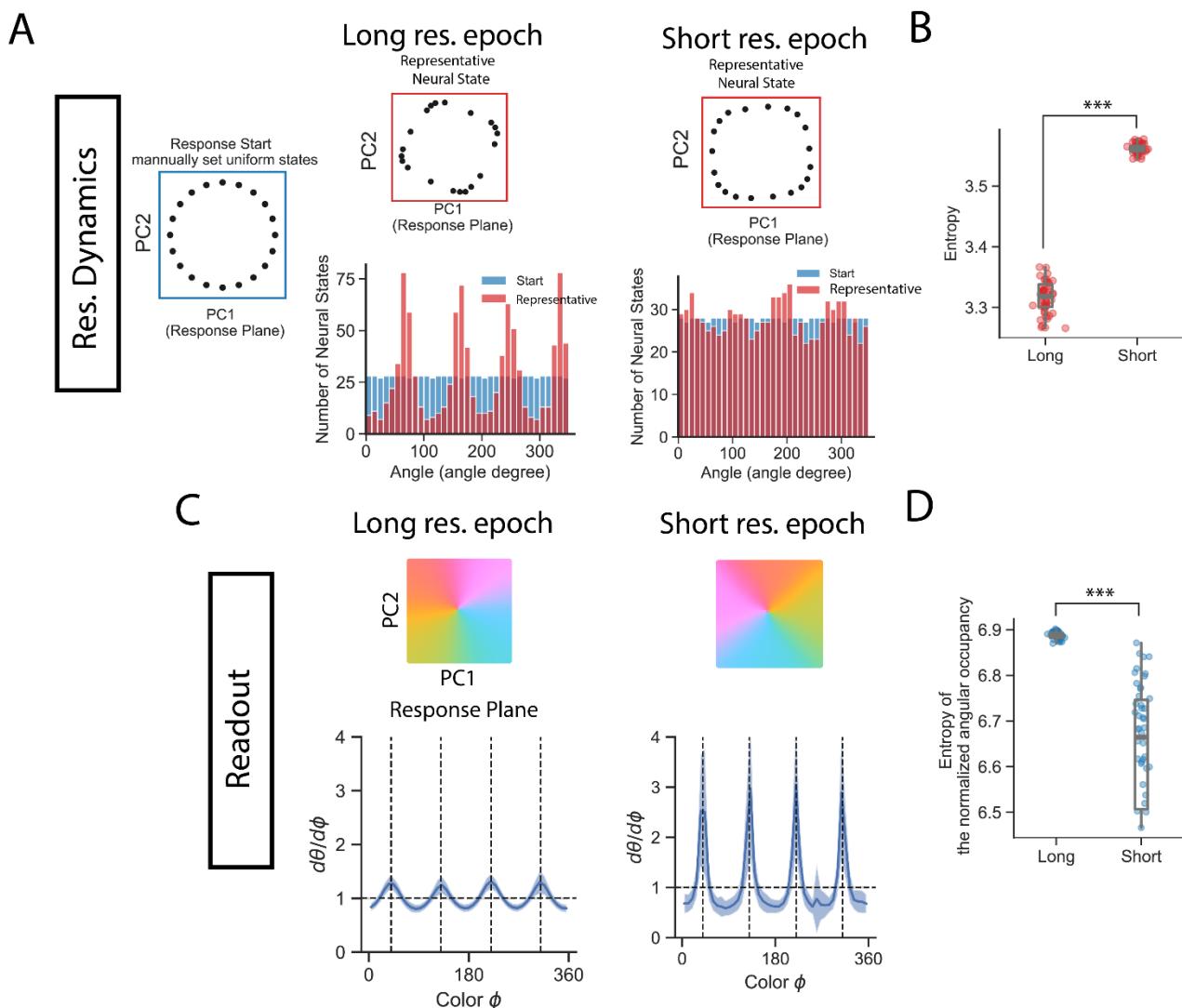
SI Figure 1. Trained Biased RNN ($\sigma_s = 25^\circ$) exhibited behavior features similar to the experiments.
(A) The distribution of memory error. Each RNN was tested on 1000 trials with uniformly distributed input color. Memory error from 50 Biased RNNs were concatenated. Bottom: experimental data. (B) Upper: the distribution of output color when the input colors were sampled uniformly for each trial. 50 Biased RNNs results were concatenated, 1000 trials each. Black dashed lines indicate common colors in the training prior distribution. Bottom: Human's output color distribution. Note human had not yet been trained on biased prior color distribution. The biasness observed here may come from the biased color prior in the natural environment. (C) Upper: memory error for each input color. 50 Biased RNN was used, each input color was run through 1000 trials for each of the Biased RNN. The error band is the standard error of the mean across trials and RNNs. The delay lengths were fixed to either 0.1s or 1s. Inset shows the training environmental prior. Bottom: human participants were trained on a biased prior distribution (inset). Their performance during the last third of the training trials is shown. Color represents results when delay was fixed to either 0.5s or 4.0s. Experimental figures were plotted using the data published by Panichello et al¹.

998



999

1000 **SI Figure 2. Eigenvalues of the fixed points.** (Left) We trained 50 RNNs for each prior distribution. For
1001 each RNN, we searched for the attractors (see Methods). The largest eigenvalue (real component) for each
1002 attractor is denoted as $v_{i,j}$, where i indicates the i th RNN and j indicates the j th attractor. The averaged
1003 largest eigenvalues across different attractors within the same RNN is denoted as $u_i = \frac{1}{N_j} \sum_j v_{i,j}$. Dots in the
1004 figure represent the mean of largest eigenvalues across different RNNs, $\frac{1}{N_i} \sum_i u_i$. Error bars indicate the
1005 standard deviation of u_i across different RNNs. (Right) Same as left but considering the largest real
1006 component of saddles. Both attractors and saddles become stronger (larger absolute values for the real parts)
1007 for narrower prior distributions.
1008



1009
1010
1011 **SI Figure 3. RNNs trained on a shorter response epoch have weaker response dynamics bias and**
1012 **stronger readout bias.** Training prior $\sigma_s = 3^\circ$. Computing details are identical to Figure 6 except that the
1013 RNN models were replaced by which trained with long/short response (res.) epoch. (A) Long/short response
1014 epoch: one example RNN trained with long (200 ms)/short (40 ms) response epoch duration. Compared to
1015 the RNN trained on long response epoch, RNN trained on short response has smaller change of neural states
1016 during the response epoch. (B) Increasing entropy of the representative neural state distribution also suggests
1017 that RNNs trained on short response epoch have less neural states' angular changes during the response.
1018 Each dot is one RNN's representative neural state entropy (red histogram in panel A). ***: $p < 10^{-3}$
1019 Wilcoxon signed-rank test. (C) RNNs trained on short response epoch have more biased readout. Line shows
1020 the mean of 50 RNNs, and error band is the standard deviation. Dash lines: four common colors. Outlier
1021 RNNs were removed. (D) Entropy of angular occupancy, dot is an entropy of one RNN's angular occupancy
1022 (normalized). Outliers were not shown. ***: $p < 10^{-3}$ Wilcoxon signed-rank test.