

1 Title:

2 Unsupervised learning for robust working memory

3

4 Authors:

5 Jintao Gu¹ and Sukbin Lim^{1,2*}

6

7 Affiliation:

8 ¹ Neural Science, New York University Shanghai, 1555 Century Avenue, Shanghai, 200122, China

9 ² NYU-ECNU Institute of Brain and Cognitive Science at NYU Shanghai, 3663 Zhongshan Road North,
10 Shanghai, 200062, China

11

12 *corresponding author:

13 Email: Sukbin.lim@nyu.edu

14 Abstract

15 Working memory is a core component of critical cognitive functions such as planning and decision-
16 making. Persistent activity that lasts long after the stimulus offset has been considered a neural
17 substrate for working memory. Attractor dynamics based on network interactions can successfully
18 reproduce such persistent activity. However, it suffers from a fine-tuning of network connectivity, in
19 particular, to form continuous attractors suggested for working memory encoding analog signals. Here,
20 we investigate whether a specific form of synaptic plasticity rules can mitigate such tuning problems in
21 two representative working memory models, namely, rate-coded and location-coded persistent activity.
22 We consider two prominent types of plasticity rules, differential plasticity targeting the slip of instant
23 neural activity and homeostatic plasticity regularizing the long-term average of activity, both of which
24 have been proposed to fine-tune the weights in an unsupervised manner. Consistent with the findings of
25 previous works, differential plasticity alone was enough to recover a graded-level persistent activity with
26 less sensitivity to learning parameters. However, for the maintenance of spatially structured persistent
27 activity, differential plasticity could recover persistent activity, but its pattern can be irregular for
28 different stimulus locations. On the other hand, homeostatic plasticity shows a robust recovery of
29 smooth spatial patterns under particular types of synaptic perturbations, such as perturbations in
30 incoming synapses onto the entire or local populations, while it was not effective against perturbations
31 in outgoing synapses from local populations. Instead, combining it with differential plasticity recovers
32 location-coded persistent activity for a broader range of perturbations, suggesting compensation
33 between two plasticity rules.

34 Author Summary

35 While external error and reward signals are essential for supervised and reinforcement learning, they
36 are not always available. For example, when an animal holds a piece of information in mind for a short
37 delay period in the absence of the original stimulus, it cannot generate an error signal by comparing its
38 memory representation with the stimulus. Thus, it might be helpful to utilize an internal signal to guide
39 learning. Here, we investigate the role of unsupervised learning for working memory maintenance,
40 which acts during the delay period without external inputs. We consider two prominent classes of
41 learning rules, namely, differential plasticity, which targets the slip of instant neural activity, and
42 homeostatic plasticity, which regularizes the long-term average of activity. The two learning rules have
43 been proposed to fine-tune the synaptic weights without external teaching signals. Here, by comparing
44 their performance under various types of network perturbations, we reveal the conditions under which
45 each rule can be effective and suggest possible synergy between them.

46 Introduction

47 Continuous attractors have been hypothesized to support brains' temporary storage and
48 integration of analog information (1–3). An attractor is an idealized stable firing pattern that persists in
49 the absence of stimuli, and integration is allowed if these attractors form a continuous manifold.
50 Theoretical models predict that neural activity should be restricted within but free to move along this
51 manifold, making stochastic fluctuation correlated among neurons, as is validated in the brainstem
52 oculomotor neural integrator (4), the entorhinal grid cell system (5), and prefrontal visuospatial selective
53 neurons (6).

54 Computationally, the performance of continuous attractors is known to be sensitive to network
55 parameters, which is termed as the “fine-tuning problem” (7,8). The slight imperfection of synaptic
56 weight asymmetry could make continuous attractors break down into a few discrete attractors or cause
57 an overall drift of activities. This raises the question of how continuous attractors could exist in the
58 brain. Noting that the model is just an idealization, earlier studies have proposed that continuous
59 attractors can be approximated by finely discretized attractors with a hysteresis of coupled bi-stable
60 units, which would make the system more robust (9,10). Recent theoretical studies suggest other
61 complementary mechanisms, including derivative feedback and short-term facilitation, with the former
62 slowing down activity decay (11,12) and the latter transiently enhancing stability (13,14).

63 These workarounds could make continuous attractors more tolerant to parameter perturbation.
64 Not mutually exclusively, long-term plasticity is believed to take part in settling a reasonable parameter
65 range. For example, the plasticity involved in the fish oculomotor integrator has been most studied.
66 Previous works have proposed either visually supervised plasticity (15–17) or self-monitoring plasticity
67 acting in the dark (18,19). These plasticity rules utilize time-derivative signals to detect slip of eye
68 position or neural activity. Note that similar mechanisms can be generalized to mediate the tuning
69 conditions of the parametric working memory encoding analog information (11,17,20). More broadly,
70 derivative-based rules have been suggested to learn temporal relationships between input and output
71 (21–23) and in reinforcement learning (24–26).

72 Another class of long-term synaptic plasticity suggested for continuous attractors is homeostatic
73 plasticity, which regularizes the excitability of neurons (27). Many models focused on the role of
74 homeostatic plasticity to prevent instability. As homeostatic plasticity tends to pull excitation down or
75 boost inhibition when network activity is higher than a reference value, the positive feedback between
76 network activity and activity-dependent plasticity can be counterbalanced (28). On the other hand,
77 Renart et al. (29) considered network storing spatial information in a spatially localized “bump” activity
78 pattern and proposed an additional role of homeostatic plasticity, that is to regularize the network
79 patterns and recover tuning condition for spatial working memory perturbed by the heterogeneity of
80 local excitability. Similarly, Pool and Mato (30) suggested that for developing orientation selectivity
81 through Hebbian learning in recurrent connections, homeostatic plasticity can enforce symmetry in
82 synaptic connections such that all orientation can be represented equally in the networks.

83 Both differential and homeostatic plasticity suggested for attractor networks are unsupervised.
84 External supervisory or reward signals are not required to recover the required tuning condition. As
85 shown previously, they can act after the offset of sensory signals and might be suitable for memory
86 tasks that typically have a long memory period without external input. However, previous works have
87 investigated the effect of differential plasticity and homeostatic plasticity partially for different types of
88 continuous attractor or under particular types of perturbations in the network parameters.

89 Therefore, we investigated whether these two forms of learning can recover persistent activity
90 in continuous attractors, which require fine-tuning conditions of network parameters. As a systematic
91 study, we considered two different types of continuous attractors, namely, rate-coded and location-
92 coded persistent memory, in a single framework, called the negative derivative feedback mechanisms
93 (11,12). First, we formally described the fine-tuning problem in a rate-coded attractor system with a
94 simpler network architecture than a location-coded attractor. We examined the effects of differential
95 plasticity and homeostatic plasticity and how recovery from perturbation in connectivity depends on the

96 learning parameters. Then we extended the scope of our investigation to a location-coded system that
97 requires spatially structured networks and investigated the recovery of tuning conditions under various
98 types of perturbations. Finally, we demonstrated that two rules could partially compensate for each
99 other when they are combined.

100 Results

101 Rate-coded persistent activity in one homogenous population

102 Before we discuss the synaptic plasticity rule that stabilizes persistent spatial patterns of
103 activity, we first consider the similar mechanism applied for a rate-coded persistent activity where the
104 persistent firing rate of memory neurons varies monotonically with the encoded signals (2). Compared
105 to location-coded memory suggested for maintaining spatial information, the rate coded one has been
106 suggested to maintain a graded level of information such as somatosensory vibration frequency (31,32).
107 Previous theoretical works proposed that recurrent circuits can maintain both types of memory based
108 on similar feedback mechanisms despite the different network architecture (12). Thus, we first gain
109 insight into how the specific form of synaptic plasticity can stabilize persistent memory in the rate
110 coding scheme, which has a simpler network structure.

111 As the rate-coded network can be built upon a spatially homogeneous structure, its dynamic
112 principle can be captured in the mean-field equations describing the network dynamics with one
113 variable (see Methods). Two representative feedback mechanisms have been proposed based on
114 recurrent network interactions, positive feedback, and negative derivative feedback, both of which is
115 described by the following equation,

$$116 \quad \frac{dr}{dt} = -r + w_{net}r - w_{der}\frac{dr}{dt} + I(t). \#(1)$$

117 In the above equation, r represents the mean firing rate of the network activity. We considered time t
118 and other time constants are unitless (normalized with the intrinsic time constant of r) for simplicity.
119 The first and last terms on the right side represent the intrinsic leakage and transient external input. The
120 second and third terms represent the feedback arising from recurrent inputs.

121 In the positive feedback models, the excessive excitatory inputs need to be tuned to cancel the
122 intrinsic leakage such that the net gain w_{net} in the second term is tuned to be one, whereas w_{der} is
123 typically zero (11). On the other hand, in the negative derivative feedback models, balanced excitatory
124 and inhibitory recurrent inputs with different kinetics generate the resistive force against memory
125 slippage similar to time-derivative activity in the third term. As its strength represented by w_{der} increases
126 while the second term remains relatively small, the effective time constant of decay of network activity
127 increases proportionally. Thus, for large negative derivative feedback, the decay of activity slows down
128 (11).

129 With a long time-constant of decay, both networks show integrator-like properties such that
130 during the stimulus presentation, it integrates the external input. After its offset, it maintains persistent
131 activity at different levels (Fig. 1A). However, any memory circuits keeping the information in continuum
132 states face a fine-tuning problem (7,8,33). Similarly, for rate-coded persistent memory, despite the
133 different tuning conditions in positive feedback models and negative-derivative feedback models, the

134 deviation from the perfect tuning leads to a gross disruption of persistent activity. For instance, a
135 reduction in the E-to-E connection mimics the effect of NMDA perturbation in memory cells shown
136 experimentally (11,34). Such a perturbation causes an imbalance between the recurrent excitation and
137 inhibition in negative derivative feedback models and leads to the rapid decay of the activity (Fig. 1B).

138 **Figure 1. Recovery of rate-coded persistent activity through differential plasticity.**

139 A: Maintenance of persistent activity through negative-derivative feedback. With balanced excitation
140 and inhibition with slower excitation, the network can maintain persistent activity at different rates
141 (solid, dotted, and dash-dotted curves represent activity for different input strengths). B: Disruption of
142 persistent activity (bottom) under the perturbation in the recurrent excitatory connections (arrow in top
143 panel). C: Schematics of differential plasticity (top) and recovery of E-I balance under differential
144 plasticity (bottom). D: Maintenance of persistent activity after the recovery of E-I balance through
145 differential plasticity.

146 **Stabilization of persistence through differential plasticity**

147 To mitigate this fine-tuning condition and to make the network resilient against perturbations,
148 several forms of synaptic plasticity have been proposed. Two prominent synaptic plasticities suggested
149 for persistent activities are homeostatic plasticity (27,29) and differential plasticity (18,19). Here, we
150 examine how each plasticity can stabilize a rate-coded persistent activity.

151 First, we consider differential synaptic plasticity where the synaptic update depends on the
152 firing rate and its time derivative of pre- and postsynaptic activities [(18); Fig. 1C]. Previous work showed
153 that such a plasticity rule updates the synaptic connection to reduce the overall derivative of network
154 activities (18). We considered the negative-derivative feedback model composed of one homogenous
155 population to understand further how the fine-tuning condition can be achieved through the differential
156 plasticity rule. We assumed that the network receives balanced recurrent excitatory and inhibitory
157 inputs with its strengths denoted as W_{exc} and W_{inh} , and excitatory inputs have slow kinetics than the
158 inhibitory inputs. If initially balanced excitation and inhibition is perturbed by the reduction in the
159 excitatory connection and excitatory connection changes according to the differential plasticity rule, the
160 dynamics of the system can be captured by the firing rate r and excitatory connection strength W_{exc} as

$$161 \frac{dr}{dt} = -r + (W_{exc} - W_{inh})r - w_{der} \frac{dr}{dt} \quad \text{where } w_{der} \text{ is proportional to } W_{inh} \text{ and the difference of the time}$$
$$\frac{dW_{exc}}{dt} = -\alpha \frac{dr}{dt} r \#(2)$$

162 constants for excitatory and inhibitory inputs feedback (Methods).

163 The steady states of the system are $r = 0$ or $dr/dt = 0$, where the latter can be achieved for
164 balanced excitation and inhibition, that is, W_{exc} becomes closer to W_{inh} for large W_{inh} . However, in
165 successive trials where each trial is composed of stimulus presentation and delay period, and assuming
166 that external input during the stimulus presentation reset r without changing W_{exc} , we found that only
167 $dr/dt = 0$, that is, the balanced tuning condition can be achieved (Fig. 1C,D). Once this tuning condition is
168 achieved, the network can maintain the graded level of persistent activities (Fig. 1D).

169 We further investigated how the recovery of a tuning condition depends on the parameters of
170 synaptic plasticity by examining the phase plane of r and W_{exc} (Fig. 2A). The learning speed α and W_{inh}
171 have similar effects of modulating the vector field along the W_{exc} -axis such that increasing W_{inh} is

172 effectively the same as decreasing α (Fig. 2B-C; Methods). In other words, stronger derivative feedback
173 requires a longer time to recover after the same percentage of perturbation, so the recovery duration
174 should scale with synaptic strength W_{inh} . On the other hand, larger perturbation leads to the initial W_{exc}
175 further away from the balanced state, making it longer to recover its tuning condition (Fig. 2D). Finally,
176 overall input strengths during the stimulus presentation determine the magnitude of r to be reset during
177 the stimulus presentation such that larger stimulus strength pushes the system in a faster speed regime
178 and makes the system faster to converge (Fig. 2E).

179 **Figure 2. Recovery dynamics dependence on learning parameters under differential plasticity.**

180 A: Phase-plane of activity r and synaptic strength of recurrent excitation W_{exc} . The black arrows
181 represent a vector field for the dynamics of r and W_{exc} , described in Eq. 2. The red curve is a trajectory
182 starting from 10% perturbation in W_{exc} , that is, $W_{exc} = 0.9W_{inh}$ with $W_{inh} = 500$. During the stimulus
183 presentation, the trajectory jumps horizontally, and input strengths vary randomly across trials. B-D:
184 Dependence of recovery speed on learning and network parameters. The minimum number of trials for
185 W_{exc} to reach up to 1% precision was obtained by varying the learning speed α (B), W_{inh} (C), perturbation
186 strength (D), and relative mean input strengths across the trial (E).

187

188 **Homeostatic plasticity is effective but sensitive**

189 While differential plasticity has been shown to stabilize the rate-coded persistent activity
190 (11,18,19), homeostatic plasticity has been suggested to stabilize different forms of memory, such as
191 spatial working memory (29) and discrete working memory (35,36). The homeostatic plasticity regulates
192 the excitability of postsynaptic neurons such that in its typical form, all incoming synapses onto the
193 postsynaptic neurons multiplicatively scale for the long-term average rate to achieve their target firing
194 rates r_0 (Fig. 3A). As for differential plasticity, we examined the effect of homeostatic plasticity in one
195 homogenous population for a rate-coded persistent activity, whose dynamics is described as

$$\begin{aligned} \frac{dr}{dt} &= -r + (W_{exc} - W_{inh})r - w_{der} \frac{dr}{dt} \\ \frac{dW_{exc}}{dt} &= -\alpha W_{exc}(r - r_0). \end{aligned} \quad (3)$$

197 The steady-state of such a system is achieved when $r=r_0$ and $dr/dt = 0$, that is, $W_{exc} \approx W_{inh}$ for large W_{inh} .
198 Note that this is more stringent than those for differential plasticity that requires the latter balance
199 condition only.

200 **Figure 3. Recovery of rate-coded persistent activity through homeostatic plasticity.**

201 A: Schematics of homeostatic plasticity scaling the strengths of incoming synapses to achieve the target
202 firing rate r_0 . B-C: Recovery of E-I balance (B) and maintenance of persistent activity at the different
203 levels after the recovery (C).

204

205 Similarly to differential plasticity, we found that the steady-state can be achieved through
206 homeostatic plasticity. However, it requires additional tuning of input strengths. The input strength set
207 the initial condition for r at the beginning of the delay and the mean of initial r needs to be r_0 on

208 average. In such a case, the network achieves the balance condition (Fig. 3B) and maintains the rate-
209 coded persistent activity (Fig. 3C). However, for inadequately tuned input, the steady-state cannot be
210 achieved (Fig. 4A-B). For the mean of the input strength making r in the beginning of the delay period
211 smaller (larger) than r_0 , the dynamics of r drifts upward (downward) to achieve r_0 on average during the
212 delay period (Fig. 4A-B, lower panels). Consequently, W_{exc} is stabilized to be excessive (deficient)
213 compared to W_{inh} (Fig. 4A-B, upper panels). Thus, the homeostatic rule for rate-coded persistent
214 memory requires tuning of input strengths and duration to achieve its target rate r_0 and balance
215 condition of recurrent excitation and inhibition.

216 **Figure 4. Sensitivity of homeostatic learning rule on learning parameters.**

217 A-B: Sensitivity to mismatch between r_0 and input strengths. With lower mean input strengths
218 compared to those in Fig. 3, the mean firing rate at the beginning of the delay period is lower than the
219 target firing rate r_0 , and the dynamics drift upward to achieve r_0 on average during the delay period (A,
220 bottom). This results in excessive W_{exc} compared to W_{inh} (A, top). The opposite leads to the decay of
221 activity and deficient W_{exc} (B). C: Sensitivity to the learning rate. For a faster learning rate, the
222 homeostatic plasticity leads to the oscillation even for property tuned inputs, and the activity can vary
223 across different trials for the same strength of the input.

224

225 Also, it is notable that the stability analysis further reveals that near the steady-state, the system
226 shows damped oscillation. Its frequency depends on the speed of the homeostatic learning rule such
227 that faster learning leads to faster oscillation. In successive trials with reset in r , the faster learning leads
228 to the ongoing oscillation near the steady-state even for a properly tuned input such that for different
229 trials, the dynamics cannot be stabilized (Fig. 4C). Overall, the analysis of one homogenous population
230 shows that although homeostatic rule can recover persistent activity for rate-coded memory, it is
231 sensitive to input parameters and learning speed.

232 **Location-coded persistent memory in spatially structured network**

233 So far, we showed how two prominent plasticity rules could stabilize rate-coded persistent
234 memory in one homogenous population. However, whether the same mechanism can be generalized to
235 stabilize location-coded persistent memory is in question because both rules are local, depending on
236 pre- and postsynaptic activity but have no regularization on a spatial pattern of activities required for
237 encoding spatial information. Here, we considered the negative derivative feedback model suggested
238 for spatial working memory (12) and explored under which condition such generalization can be made.

239 Previous work showed that the principle for negative derivative feedback found for one
240 homogenous population could be extended to the network with a functionally columnar structure
241 required to maintain a spatial pattern of persistent activity. Consistent with experimental observations
242 (37–39), both excitatory and inhibitory neurons in each column have similar spatial selectivity. The
243 connectivity strengths decrease as the preferred features over the columns get dissimilar (Fig. 5A-B).
244 Assuming translation-invariance of connectivity strength such that it depends only on the distance
245 between neurons' preferred features, the network activity is symmetric under the translation of
246 stimulus location.

247 **Figure 5. Location-coded persistent activity and its disruption under perturbation of tuning**

248 A: Schematics of the spatial structure of network for location-coded memory. We considered that both
249 excitatory and inhibitory neurons are organized in a columnar structure where each column consists of
250 neurons with a similar preferred feature of the stimulus. Blue and red represent excitatory and
251 inhibitory connections, respectively. B: Example connectivity matrix showing symmetry under
252 translation. We considered the stimulus feature neurons encode the spatial information during the
253 delay period, which lies on a circle, represented by θ ranging between $-\pi$ and π . We assumed
254 translation-invariance with the synaptic strengths depending only on the difference between preferred
255 features of post and presynaptic neurons. C: Decomposition of spatially patterned activity into Fourier
256 modes. Under translation-invariance, the activity can be decomposed into Fourier modes, and with
257 strong negative derivative feedback, the dynamics of each mode become independent. D: Location-
258 coded persistent activity under E-I balance. The activity during five consecutive trials was shown where
259 the center of input is shifted randomly, showing maintenance of the spatial pattern of activity (upper
260 panel) as well as elevated persistent activity at the stimulation center (lower panel). E: Disruption of
261 persistent activity under 10% global perturbation in the E-to-E connection.

262

263 Under translation-invariant connectivity and activity patterns, dynamics can be analyzed
264 through Fourier analysis ((12); Fig. 5C). For a large recurrent excitation and inhibition, the dynamics of
265 Fourier modes are approximately independent of each other, each of which is analogous to the
266 dynamics of one homogeneous population. Thus, the condition for negative derivative feedback in each
267 Fourier mode is similar to the rate-coded network - slower recurrent excitation with the same condition
268 on the synaptic time constants as in the homogeneous case, and balanced recurrent excitation and
269 inhibition of that mode represented in terms of the Fourier coefficients of the synaptic strengths.

270 With similar balanced tuning conditions for the location-coded persistent memory, the
271 perturbation to the synaptic connections leads to a similar disruption in the activity as in the rate-coded
272 network (Fig. 5D-E). We first considered the multiplicative scaling down of all E-to-E connections, called
273 a global perturbation (Fig. 5E). This leads to imbalanced excitation and inhibition and decay of activity in
274 all Fourier modes. Note that the translation-invariant property is maintained under the global
275 perturbation of the connectivity. Thus, the activity pattern is still symmetric for different stimulus
276 locations despite its rapid decay to the baseline compared to the unperturbed case.

277 **Effects of plasticity under global perturbation**

278 Next, we examine whether differential plasticity and homeostatic plasticity can recover the
279 balance tuning condition for a spatially structured network. We assumed that the stimulus location
280 across different trials is uniformly distributed and changes fast enough compared to the speed of the
281 synaptic plasticity. Even if the connectivity is translation-invariant initially, stimulus at a particular
282 location can lead to asymmetrical updates in the synaptic connections. Such asymmetrical updates can
283 be mediated by slow learning and random stimulus locations having uniform distribution (30).

284 Under the maintenance of translation-invariance, the differential rule was shown to recover
285 persistent activity having spatial patterns (Fig. 6). Unimodal activity peaked at the stimulus location can
286 be maintained at any location after the differential plasticity rule recovers the balance of excitation and
287 inhibition (Fig. 6A-B). We quantified the ability to maintain location-coded persistent memory using the
288 decoding accuracy of spatial information at the end of the delay period (Methods). Initially, after global

289 perturbation, the decoding error became around one, indicating loss of spatial selectivity, but over the
290 course of learning with differential plasticity, it becomes close to zero (Fig. 6C). In line with this, the time
291 constant of decay of different Fourier modes was shown to prolong (Fig. S1). In the eigenvector
292 decomposition of the connectivity matrix, eigenvectors corresponding to the leading eigenvalues were
293 found to be similar to Fourier modes, which is a signature of preservation of translation-invariance ((40);
294 Fig. S1). The ratios of associated eigenvalues increase to one, albeit the different speeds, suggesting the
295 recovery of balance tuning condition in each mode (Fig. 6D).

296 **Figure 6. The effect of differential plasticity under weak global perturbation.**

297 A: Recovery of location-coded memory with differential plasticity under 10% global perturbation in the
298 E-to-E connections. B: Activity pattern at the end of delay period after the recovery. With the
299 connectivity frozen at trial 8000 (arrow in C), the spatial pattern of activity at the end of the delay period
300 was shown for different stimulus locations. C: Improvement of decoding accuracy with learning. An
301 individual trial refers to one memory task with a specific stimulus location. For each trial, we took the
302 snapshot of activity as in B and quantified the decoding error using the population vector decoder
303 (Methods). Dashed line indicates decoding error before perturbation D: Recovery of E-I balance for
304 different Fourier modes. The eigenvector decomposition reveals the effective time constant of decay
305 and recovery of E-I balance in different Fourier modes (Methods; Fig. S1). E: Mean (black) and standard
306 deviation (red) of spatial selectivity across neurons quantified by the first Fourier component of each
307 neuron's tuning curve at the end of the delay period. F: Normalized standard deviation of Fourier tuning
308 in (E) where its decrease with learning indicates recovery of translation invariance.

309 However, if translation-invariance breaks down, then Fourier analysis cannot be applied. This
310 breakdown can occur either when the learning is too fast such that it cannot overcome asymmetry
311 introduced by the different stimulus location at each trial or when the perturbation is too strong such
312 that the activity of some neurons is stabilized to zero. Figure 7 shows the latter case – for stronger
313 perturbation, the persistence of activity is recovered under differential plasticity, but the spatial pattern
314 is fragmented by silent neurons (Fig. 7A-B). The decoding accuracy still improves during learning due to
315 active neurons encoding spatial information (Fig. 7C and 7E).

316 **Figure 7. The effect of differential plasticity under larger global perturbation.**

317 A-B: Time course and tuning of activity after the recovery of E-I balance under differential plasticity
318 (arrow in C), but 15% global perturbation was used. Differential plasticity stabilizes some neurons'
319 activity to zero, making them unresponsive to all stimulus locations. C: Improvement of decoding
320 accuracy despite silent neurons due to compensation by neighboring active neurons. D: Normalized
321 standard deviation of the first Fourier component over neurons, which remains high at the end of trials.
322 E: Decoding errors under different strengths of perturbations. Dashed line indicates decoding error
323 before perturbation. F: mean (black) and standard deviation (red) of spatial selectivity (F) under
324 different strengths of perturbations.

325

326 For larger perturbation, the recovery of persistent activity and decoding accuracy is not uniform
327 across different neurons as translation-invariance breaks down. To quantify this heterogeneity, we
328 calculated the first Fourier component of the tuning curve of each neuron at the end of the delay

329 period, representing its spatial selectivity, and obtained its mean and variance across neurons (Fig. 6E,
330 7D, and 7F; Method). Its mean increases with learning, indicating the increase of spatial selectivity with
331 learning, shown for a broader range of perturbation (Fig. 6E and 7F). For a relatively weak global
332 perturbation, the variance can transiently increase, reflecting an overall increase of activity level, but the
333 normalized variance by the mean decreases for smaller perturbation with the translation-invariance
334 maintained (Fig. 6E-F). However, for larger perturbation, such normalized variance is not reduced to
335 zero even after decoding accuracy reaches its asymptote, indicating the breakdown of translation
336 invariance (Fig. 7D and 7F).

337 While the maintenance of translation-invariance is not guaranteed under differential plasticity,
338 homeostatic plasticity has been suggested to recover translation-invariance perturbed under cellular
339 heterogeneity or other types of synaptic plasticity such as Hebbian learning (29). Indeed, the application
340 of homeostatic learning rule to negative derivative feedback network recovers persistent unimodal
341 activity at different locations (Fig. 8). As for differential plasticity, decoding accuracy improves with
342 learning as the E-I balance is recovered (Fig. 8C and 8E). In contrast to differential plasticity, homeostatic
343 plasticity achieved a low variance of spatial selectivity across neurons for a broader range of global
344 perturbation, suggesting the maintenance of translation-invariance (Fig. 8D and F). As for a rate-coded
345 memory network, homeostatic plasticity requires the input strengths to be tuned to match average
346 delay activity to the target rate. However, the condition on the input strength can be mitigated with
347 cellular or synaptic nonlinearity for location-coded persistent memory. The information can be decoded
348 from the peak of the bump activity and is not sensitive to the amplitude of the bump.

349 **Figure 8. The effect of homeostatic plasticity under global perturbation.**

350 A-D: Same as Fig. 7A-D but with homeostatic plasticity, showing recovery of location-coded persistent
351 activity and translation-invariance. E: Recovery of E-I balance in different Fourier modes shows the
352 same recovery speed because homeostatic plasticity multiplicatively scales all afferent weights. F:
353 Decoding errors (black) and normalized deviations of spatial selectivity (red) for different perturbation
354 strengths. Dashed line indicates decoding error before perturbation.

355

356 **Effects of plasticity under local perturbation**

357 We further investigated the effect of differential and homeostatic plasticity, where the balance
358 of excitation and inhibition is locally perturbed. We considered two different types of local perturbations
359 – first, postsynaptic perturbations, where synaptic strengths projected onto a particular group of
360 neurons were perturbed (Fig. 9). For instance, this can be incurred by perturbation in NMDA receptors,
361 which is considered to be prominent in the E-to-E connections (41). Mathematically, it is analogous to a
362 row-wise perturbation in the E-E connectivity matrix (Fig. 9A). Another type of perturbation is the
363 presynaptic one, where outgoing synaptic strengths are perturbed (Fig. 10). This perturbation can be
364 caused by reducing transmitter release and is analogous to column-wise perturbation in the connectivity
365 matrix (Fig. 10A). We considered a smooth bell-shaped perturbation assuming that the neurons with
366 similar preferred spatial selectivity are clustered, and the effect of local perturbation dissipates across
367 the clusters (42).

368 **Figure 9. The effect of differential and homeostatic plasticity under postsynaptic perturbations.**

369 A: Schematics of postsynaptic perturbations where the rows of the connectivity matrix are multiplied by
370 different scaling factors. Perturbation is centered at $\theta=0$ and bell-shaped. B: Activity pattern under 15%
371 post-synaptic perturbations before any plasticity. C-D: activity pattern shaped by the differential (C) and
372 homeostatic (D) plasticity. E-F: Decoding errors (black) and normalized deviations of spatial selectivity
373 (red) for different perturbation strengths after applying differential (E) and homeostatic (F) plasticity.
374 Under differential plasticity, some neurons were silenced near the perturbation site (C), and the
375 translation-invariance breaks down albeit with decent decoding performance (E). In contrast,
376 homeostatic plasticity recovers the location-coded persistent activity for a broad range of postsynaptic
377 perturbations (F).

378 **Figure 10. The effect of differential and homeostatic plasticity under pre-synaptic perturbations.**

379 A: Schematics of pre-synaptic perturbations where the columns of the connectivity matrix are multiplied
380 by different scaling factors. B-F: Same as in Fig. 9B-F but under 15% pre-synaptic perturbation. Unlike
381 postsynaptic perturbations, differential plasticity recovers persistent activity and translation-invariance
382 (C, E). In contrast, with homeostatic plasticity, the activity pattern was distorted, resulting in worse
383 decoding accuracy and translation-invariance for larger perturbations (D, F).

384

385 We first examined the effect of plasticity in postsynaptic perturbations. In negative derivative
386 feedback models, the postsynaptic perturbation disrupts local E-I balance, leading to quick decay of
387 activity in the vicinity of the perturbed site (Fig. 9B). Under relatively weak perturbation, both
388 differential and homeostatic plasticity can recover E-I balance and the ability to maintain persistent
389 activity at the perturbed site (Fig. 9E-F). However, when the perturbation becomes larger, differential
390 and homeostatic plasticity show different recovery patterns as for the global perturbation (Fig. 9C-D).
391 While homeostatic plasticity recovers both persistent activity and translation-invariance, differential
392 plasticity persistently silences some neurons and cannot recover translation invariance. The fragmented
393 spatial activity results in a high variance of spatial selectivity, breaking down translation-invariance,
394 while the decoding accuracy is still good with compensation by higher activity at the vicinity of silent
395 neurons (Fig. 9E). In contrast, homeostatic plasticity efficiently recovers translation-invariance caused by
396 the overall reduction of synaptic strengths onto particular neurons as it multiplicatively scales those
397 connections (Fig. 9F).

398 Next, we considered the effect of plasticity under presynaptic perturbations, which showed
399 better performance of differential plasticity than homeostatic plasticity (Fig. 10). As in the postsynaptic
400 perturbations, presynaptic perturbation causes activity at the perturbed site to decay because
401 perturbation in outgoing synapses mostly affects the incoming synapses of neurons with similar spatial
402 selectivity (Fig. 10B). Differential plasticity can recover persistent activity and translation-invariance for a
403 broad range of presynaptic perturbation (Fig. 10C, 10E). On the other hand, homeostatic plasticity
404 cannot stabilize persistent activity for relatively large perturbation, and the distortion of activity pattern
405 is more substantial near the perturbed site (Fig. 10D). This is because presynaptic perturbation
406 introduces an asymmetry in the synaptic strengths projecting onto neurons near the perturbed sites,
407 which cannot be recovered through homeostatic plasticity that regulates the overall scaling of incoming
408 synapses. Thus, although the average postsynaptic activity is recovered through increased excitability,
409 the bump activity drifts towards instead of away from the perturbed site after learning, leading to a low
410 decoding accuracy and breakdown of translation-invariance both (Fig. 10F).

411 **Effect of combining differential and homeostatic plasticity**

412 As differential plasticity and homeostatic plasticity are effective in recovering persistent activity
413 and translation-invariance under the different types of perturbations, we examined whether the
414 combination of these two plasticities can utilize the advantage of both plasticities. Following the
415 previous models considering the combination of Hebbian and homeostatic plasticity, we considered a
416 multiplicative combination of two plasticities where differential plasticity replaces the Hebbian learning.
417 The synaptic connection from neuron j to neuron i is expressed as a product of two variables, $W_{ij} = g_i U_{ij}$
418 with the dynamics of g_i and U_{ij} are given as

$$419 \quad \begin{aligned} \frac{dg_i}{dt} &= -\alpha_h(r_i - r_0)g_i \\ \frac{dU_{ij}}{dt} &= -\alpha_d \frac{dr_i}{dt} r_j. \#(4) \end{aligned}$$

420 In the above equations, g_i reflects the homeostatic scaling, and U_{ij} evolves according to differential
421 plasticity, with the learning rates given as α_h and α_d , respectively.

422 We first examined the effect of combined plasticity under global perturbations. Although
423 differential plasticity alone can lead to the silence of activity, homeostatic plasticity prevents it by
424 boosting lower-than-target activity. Thus, combined plasticity could recover location-coded persistent
425 activity and translation invariance for a broad range of perturbations (Fig. 11A and 11D). Note that such
426 a recovery is sensitive to the learning rates of the plasticity such that the overall speed of both
427 differential and homeostatic plasticity needs to be slow, but homeostatic one needs to be relatively fast.
428 This is because too fast homeostatic learning leads to oscillation, yet it needs to be fast enough to
429 prevent the breakdown of translation invariance introduced by differential plasticity.

430 **Figure 11. The effect of the combination of differential and homeostatic plasticity.**

431 A-C: Recovery of location-coded persistent activity under combined plasticity after 15% global (A),
432 postsynaptic (B), and pre-synaptic perturbation (C). Activity pattern after 30% local perturbations is
433 shown in Fig. S2. D-F: Decoding errors (black) and normalized deviation of spatial selectivity (red) for
434 different perturbation strengths.

435

436 The combined plasticity also shows the compensation under both types of local perturbations
437 (Fig. 11B-C, E-F). For a postsynaptic perturbation under which differential plasticity could not recover the
438 translation-invariance, the combined one shows the extension of the recovery (Fig. 9C vs. 11B, 9E vs.
439 11E). The superiority of the combined one is similar for a presynaptic perturbation under which
440 homeostatic plasticity could not recover both persistent activity and translation-invariance (Fig. 10D vs.
441 11C, 10F vs. 11F). Note that still, for a larger perturbation, the activity pattern can be distorted, and
442 translation-invariance is not perfectly recovered (Fig. 11E-F, S2). However, the decoding accuracy is
443 decent for a broad range for the combined plasticity.

444 Discussion

445 In this work, we investigated the effects of local and unsupervised learning on the stabilization of
446 persistent activity in two representative working memory models encoding analog values, namely, rate-
447 coded and location-coded persistent memory. We examined the effects of differential plasticity and
448 homeostatic plasticity by systematically varying the magnitude and form of perturbations in synaptic
449 connections. Consistent with the findings of previous works, differential plasticity alone was enough to
450 recover a graded-level persistent activity in a homogeneous population (11,18). On the other hand,
451 recovery by homeostatic plasticity requires the tuning of learning parameters. For the maintenance of
452 spatially structured persistent activity, differential plasticity could recover persistent activity, but its
453 pattern can be irregular for different stimulus locations. On the other hand, homeostatic plasticity
454 shows robust recovery of translation-invariance against particular types of synaptic perturbations, such
455 as perturbations in incoming synapses onto the entire or local populations, which are similar to the
456 inhomogeneity of neuronal gain considered previously (29). However, homeostatic plasticity was not
457 effective against perturbations in outgoing synapses from local populations. Instead, combining it with
458 differential plasticity recovers the location-coded persistent activity for a broader range of
459 perturbations.

460 Persistent activity sustained in the absence of external stimuli has been suggested as a signature
461 of static attractor dynamics. While most attractor models are based on positive feedback mechanisms,
462 we considered negative derivative feedback mechanisms with two advantages to investigate the effect
463 of synaptic plasticity. First, a negative derivative feedback network has similar tuning conditions in both
464 the rate-coded and location-coded persistent memory with the same condition on the balanced
465 excitation and inhibition and additional symmetry of translation-invariance for location-coded memory
466 (12,33). Thus, the analysis of the effect of synaptic plasticity in a relatively simple rate-coded network
467 could be extended to that in a location-coded network. Second, the negative derivative feedback model
468 is less dependent on a specific form of intrinsic nonlinearity of neurons, so graded perturbation causes a
469 graded change in the network's behavior. On the other hand, intrinsic nonlinearity plays a critical role in
470 positive feedback networks, which leads to additional complexity in systematically investigating the
471 effect of synaptic plasticity (43) (but see (44,45)). However, note that our main findings may still be valid
472 regardless of the specific underlying mechanism of working memory. For rate-coded persistent activity,
473 previous works have shown that differential plasticity recovers the tuning condition of the positive
474 feedback mechanism (17,18). We tested two plasticity rules on the location-coded memory based on
475 positive feedback and verified that homeostatic plasticity is effective against post- but not presynaptic
476 perturbation; however, differential plasticity can effectively stop drift but may lead to an irregular
477 pattern, and combination of both provides a partial remedy (Fig. S3).

478 Stable memory formation under the mixture of different forms of synaptic plasticity has been
479 proposed previously, mainly for discrete attractor networks (35,36,46). In these studies, Hebbian
480 synaptic plasticity has been suggested to form auto-associative memory guided by external inputs. To
481 prevent instability caused by Hebbian learning, compensatory mechanisms, such as homeostasis or
482 short-term plasticity, were required, which must act on a timescale similar to that of Hebbian learning
483 (47). Our work also suggests synergistic interplay between different types of plasticity, differential, and
484 homeostatic plasticity, in particular for stabilizing location-coded persistent memory. However, we note
485 that differential plasticity alone is stable. The role of homeostatic plasticity is to support translation-
486 invariance in a ring-like architecture of recurrent connections (29,30). Thus, the fast dynamics of

487 homeostatic plasticity are not required, and excessively fast dynamics can be detrimental due to
488 oscillatory instability. The interplay between anti-Hebbian learning and activity-dependent synaptic
489 scaling has been proposed for rate-coded persistent memory (48), where the anti-Hebbian rule itself
490 stabilizes the network activity and no fast homeostasis is required, as in our work.

491 In this work, we assumed the existence of synaptic plasticity only during the delay period.
492 However, differential plasticity might make the network “unlearn” if it operates the same way during
493 the stimulus period as in the delay period because the activity rise during that time would be
494 interpreted as positive drift by the plasticity. It is thus essential to constrain derivative-driven learning
495 only during a period in which the activity ought to be stabilized. In oculomotor literature, such as (17),
496 this is done by filtering fast-changing activity that is potentially related to burst of the saccadic signal.
497 When we consider the biological implementation of differential plasticity, there is an alternative way
498 that this can happen. Xie and Seung (2000), motivated by (23), showed that spike-timing-dependent
499 plasticity (STDP) is intrinsically sensitive to the time derivative of activities, and it can be approximated
500 by the differential plasticity considered in our work when the overall potentiation and depression are
501 balanced (18). Alternatively, Nygren et al. (2019) showed that similar differential plasticity could be
502 implemented through cancelation with a delayed feedback signal analogous to the derivative feedback
503 (19). In both cases, the derivative is approximated for slowly changing neural activity, and higher
504 frequency changes are filtered out. On the other hand, continuous learning with homeostatic plasticity
505 may require the adjustment of learning parameters because the long-term average firing rates of
506 neurons must reflect activity during the entire session.

507 Constraining activity drifts of individual neurons might require stricter conditions than what is
508 required to achieve stable coding of information during the memory period. While traditional
509 experimental work identified memory neurons that showed elevated persistence firing with stimulus
510 selectivity (42), recent population-level analysis revealed the stable readout of information across
511 various time points despite the diverse temporal dynamics of individual neurons (49,50). Such dynamic
512 activity in individual neurons may reflect activity in the downstream population that combines stimulus-
513 encoding persistent activity and time-varying activity, possibly reflecting time information (51,52). On
514 the other hand, memory networks themselves can allow time-varying activity such that attractor
515 dynamics are formed along with the particular activity pattern or mode, while allowing temporal
516 fluctuation along with other modes (50,53). For the latter, synaptic plasticity based on the global error
517 signal has been suggested, which can be a self-supervised signal, such as a drift in the readout activity
518 (20) or a difference from the target signal (54). Note that the resulting form of synaptic plasticity is
519 similar to differential plasticity, where the activity drift of individual neurons in differential plasticity is
520 replaced with the global error signal. Homeostatic processes, such as intrinsic plasticity, inhibitory
521 plasticity, and synaptic scaling, have also been proposed to elongate memory traces in the presence of
522 dynamic activity (48,55). In these works, the memory is maintained by a network with minimally
523 structured connectivity, and the sensitivity to learning parameters has not been analyzed.

524 Overall, our work demonstrates how unsupervised learning can mediate fine-tuning conditions for
525 working memory implementing continuous attractors. It aligns with previous works emphasizing the role
526 of unsupervised learning to generate a basis of activity patterns and dynamics underlying cognitive
527 functions (56–58). While we focused on unsupervised learning rules regularizing temporal patterns in
528 the absence of input, they can be combined with other learning rules that can act under the guidance of
529 external inputs and may make memory networks robust for a broader range of perturbations. Also, we

530 considered perturbation and synaptic plasticity only in a specific connection, recurrent E-to-E
 531 connections, but the plasticity of other connections, such as inhibitory plasticity (59–61), has been
 532 suggested to tune network homeostasis and EI balance. Given the importance of balance and
 533 homeostasis in memory circuits, further investigation is needed to examine the effect of unsupervised
 534 plasticity on various synapses. Also, to understand how the learning parameters of these plasticity rules
 535 match with neural activity, a detailed investigation of the underlying biophysical mechanisms needs to
 536 be done, possibly in models involving multiple subcellular compartments.

537 Methods

538 All codes are available at https://github.com/jtg374/NDF_ringNet_plasticity

539 Simple rate model for a homogeneous population

540 In this section, we show the derivation of a one-dimensional differential equation in Eq. 1 (see
 541 more biological structure and conditions in Lim & Goldman, 2013). For this, we considered one
 542 homogeneous population receiving recurrent excitation and inhibition with different kinetics, described
 543 by three-dimensional differential equations

$$544 \quad \tau \frac{dr}{dt} = -r + w_{exc}s_{exc} - w_{inh}s_{inh} + I(t) \#(5)$$

545 , where three dynamic variables are firing rate r , recurrent excitatory currents s_{exc} , and recurrent
 546 inhibitory currents s_{inh} . We assumed that s_{exc} and s_{inh} are low-pass filtered r with time constants τ_{exc} and
 547 τ_{inh} , respectively, and $s_{exc} \approx r$ when r hardly changes.

548 With $s_{exc} - s_{inh} \approx -(\tau_{exc} - \tau_{inh})dr/dt$, the above equation can be approximated as a one-dimensional
 549 differential equation, given as

$$550 \quad \begin{aligned} \tau \frac{dr}{dt} &= -r + (w_{exc} - w_{inh})s_{exc} + w_{inh}(s_{exc} - s_{inh}) + I(t) \\ &\approx -r + (w_{exc} - w_{inh})r - w_{inh}(\tau_{exc} - \tau_{inh})\frac{dr}{dt} + I(t) \#(6) \end{aligned}$$

551 With $w_{exc} - w_{inh}$ and $w_{inh}(\tau_{exc} - \tau_{inh})$ denoted by w_{net} and w_{der} , Eq. 6 is the same as Eq. 1. Such one-
 552 dimensional approximation allows analytic investigation on the effects of differential plasticity and
 553 homeostatic plasticity in Eq. 2 and Eq. 3.

554 In Eq. 5, when we normalized w_{exc} with w_{inh} denoted as $w = w_{exc}/w_{inh}$ and assumed w_{inh} is large
 555 such that $w_{der} \gg \tau$, Eq. 6 becomes

$$556 \quad (\tau_{exc} - \tau_{inh})\frac{dr}{dt} = (w - 1)r + \frac{I(t)}{w_{inh}} \#(7)$$

557 During the delay period with no external input $I(t)$, the dynamics with the differential plasticity in Eq. 2
 558 becomes

559
$$\begin{aligned} (\tau_{exc} - \tau_{inh}) \frac{dr}{dt} &= (w - 1)r \\ \frac{dw}{dt} &= -\frac{\alpha}{w_{inh}} \frac{dr}{dt} r \#(8) \end{aligned}$$

560 Thus, varying w_{inh} has the same effect as changing the learning speed α as illustrated in Fig. 2B,C. On the
561 other hand, the dynamics with homeostatic plasticity in Eq. 3 becomes

562
$$\begin{aligned} (\tau_{exc} - \tau_{inh}) \frac{dr}{dt} &= (w - 1)r \\ \frac{dw}{dt} &= -\alpha w (r - r_0) \#(9) \end{aligned}$$

563 That is, changing w_{inh} has no effect on the recovery speed for homeostatic plasticity.

564 In Equations 1-3, we set τ and $\tau_{exc} - \tau_{inh}$ to be unit time constant 1, and initial w_{exc} , w_{inh} and w_{der}
565 are set to be 500. $I(t)$ is a step function, giving the input for 50 unit time with its strength randomly
566 distributed as 0 and 1000 during the learning, and three representative traces of $r(t)$ for $I(t) = 250, 500$
567 and 1000 were shown before and after learning. For the differential plasticity, the learning speed α is
568 0.01. For homeostatic plasticity, α is 2×10^{-5} in Fig. 3 and Fig. 4a and b and 0.001 in Fig. 4c. r_0 is 50 in Fig. 3
569 and Fig. 4c, and 80 and 25 in Fig. 4a, and b.

570 Spatial structured network model for location-coded persistent activity

571 Following Lim & Goldman 2014, we considered a network organized in a columnar architecture
572 for spatial working memory with the equations describing the dynamics given as

573
$$\tau_E \frac{d}{dt} r_E(\theta) = -r_E(\theta) + q \left(\int_{-\pi}^{\pi} \vec{W}_{EE}(\theta, \theta') s_{EE}(\theta') d\theta' - \int_{-\pi}^{\pi} \vec{W}_{EI}(\theta, \theta') s_{EI}(\theta') d\theta' + I_s(\theta - s) I_t(t) \right)$$

574
$$\tau_I \frac{d}{dt} r_I(\theta) = -r_I(\theta) + q \left(\int_{-\pi}^{\pi} \vec{W}_{IE}(\theta, \theta') s_{IE}(\theta') d\theta' - \int_{-\pi}^{\pi} \vec{W}_{II}(\theta, \theta') s_{II}(\theta') d\theta' \right) \#(10)$$

575 where subscripts E and I represent excitatory and inhibitory populations, respectively. The activity and
576 the connectivity were indexed by their preferred spatial feature, θ , ranging between $[-\pi, \pi]$. τ_E and τ_I are
577 the time constants and $q(\cdot)$ is the input-output transfer function, which is the rectified linear function
578 given as $q(x) = x$ for $x > 0$ and otherwise, 0. $\leftrightarrow W_{ij}$ ($i, j = E$ or I) is the synaptic weight and before
579 perturbation, it was taken to be translation-invariant and Gaussian-shaped as

580
$$W_{ij}(\theta, \theta') = J_{ijexp} \left(-(\theta - \theta') / \sigma_{ij}^2 \right) \#(11)$$

581 As in the homogeneous case, \vec{s}_{ij} ($i, j = E$ or I) represents the synaptic variables whose dynamics
582 is given as

583
$$\tau_{ij} s_{ij}(\theta) = -s_{ij}(\theta) + r_i(\theta) \#(12)$$

584 Importantly, the excitatory-to-excitatory (E-to-E) time constant needs to be much larger than other
585 synapses' to make derivative feedback happen (Lim & Goldman 2013). Detailed parameters used in the
586 simulation will be given in the section "Table of parameters".

587 $I_s(\theta-s)$ and $I_t(t)$ represent the spatial and temporal profiles of external stimulus where s is the
588 center of the stimulus location. $I_s(\theta)$ also has a Gaussian shape as

$$589 \quad I_s(\theta) = J_o \exp(-(\theta/\sigma_o)^2) + h_o \#(13)$$

590 $I_t(t)$ is a pulse function smoothed by a low-pass filter with time constant τ_o

$$591 \quad I_t(x) = \begin{cases} -1 - \exp(-t/\tau_o), & \text{if } t < t_{stim} \\ I_t(t_d) \exp(-(t-t_d)/\tau_o), & \text{if } t_{stim} \leq t < t_{total} \end{cases} \#(14)$$

592 Perturbation and plasticity model

593 We considered three types of perturbations in the E-to-E connections. For the global
594 perturbation, $\leftrightarrow W_{EE}$ was set to be

$$595 \quad \vec{W}_{EE,perturbed}(\theta, \theta') = p_{uniform} \vec{W}_{EE,0}(\theta, \theta') \#(15)$$

596 Postsynaptic perturbation corresponds to a row-wise change as

$$597 \quad \vec{W}_{EE,perturbed}(\theta, \theta') = p_{pre-syn}(\theta') \vec{W}_{EE,0}(\theta, \theta') \#(16)$$

598 Similarly, presynaptic perturbation corresponds to a column-wise change as

$$599 \quad \vec{W}_{EE,perturbed}(\theta, \theta') = p_{post-syn}(\theta) \vec{W}_{EE,0}(\theta, \theta') \#(17)$$

600 Where $p(\theta)$ is a smooth function of θ , given as a Gaussian function

$$601 \quad p(\theta) = 1 - p \exp(-(\theta/\sigma_p)^2) \#(18)$$

602 To recover the persistent activity, we considered two types of plasticity, differential and
603 homeostatic plasticity, described as

$$604 \quad \frac{dW_{ij}}{dt} = -\alpha_d \frac{dr_i}{dt} r_j \#(19)$$

$$605 \quad \frac{dW_{ij}}{dt} = -\alpha_h W_{ij} (r_i - r_0) \#(20)$$

606 , where α_d and α_h represents the learning rate of differential and homeostatic plasticity. In the combined
607 one in Eq. 4, $\leftrightarrow W_{ij}$ in Eq. 2 and Eq. 3 are replaced by U_{ij} and g_i , respectively. The plasticity is only applied
608 in the delay period, and to minimize the effect of the residual stimulus, we also gated the plasticity with
609 a factor $1-I_t(t)$, though it does not make much difference if we don't add it.

610 Quantify E-I balance through eigenvalue decomposition

611 In Figure 6D and 8E, we quantify the recovery of EI balance by taking the eigenvalues of the
612 weight matrices. When translation-invariance is preserved, the values of both E-to-E matrix and other
613 weight matrices will approximately be the Fourier components of the matrices and the tuning conditions
614 for n -th Fourier modes become

$$615 \quad \lambda_{EE}(n) \lambda_{II}(n) = \lambda_{EI}(n) \lambda_{IE}(n) \#(21)$$

616 where $\lambda_{ij}(n)$ is the n -th eigenvalue of $\leftrightarrow W_{ij}$. In Fig. 6D and 8E, we did the eigenvector decomposition of
 617 the weight matrix $\leftrightarrow W_{EE}$ and found the eigenvectors resembles Fourier modes and calculate E-I balance
 618 ratio in each mode from the corresponding eigenvalues.

619 Decoding error

620 We quantified the network's memory performance by decoding the stimulus at the end of the
 621 delay. Because we used a deterministic simulation, we modeled the noise post-hoc with Poisson random
 622 number generator, assuming the spike generation is random and independent across neurons. For each
 623 neuron, we multiplied its firing rate (in Hz) by 0.2 and used the product as the mean of the Poisson
 624 random number to model its spike count in 200ms. We then decoded the location from the simulated
 625 spike count n_θ with a simple population-vector decoder:

$$626 \quad \hat{s} = \text{angle} \left(\int_{-\pi}^{\pi} e^{i\theta} n_\theta(s) d\theta \right) = \text{angle} \left(\int_{-\pi}^{\pi} \cos(s) n_\theta(s) ds + i \int_{-\pi}^{\pi} \sin(s) n_\theta(s) ds \right) \#(22)$$

627 The error is quantified by the cosine distance between the decoded location and true stimulus:

$$628 \quad \text{error} = 1 - \cos(s - \hat{s}) \#(23)$$

629 For each stimulus, the random generation of spike counts was repeated 20 times and averaged.
 630 We quantified the average error across all stimulus locations by freezing the network connectivity at
 631 each trial. After perturbation, when there is no spatial selectivity at the end of the delay, \hat{s} would be
 632 uniformly distributed, and the average error would be one, while if the spatially patterned activity is
 633 persistent with no drift, the decoding error would be close to zero. In Figure S3, we divided the stimulus
 634 locations into eight groups and visualized the average error within groups to emphasize that local
 635 perturbation affects decoding accuracy differently depending on stimulus locations.

636 Spatial selectivity and translation invariance

637 The spatial selectivity of each neuron was quantified by calculating the first Fourier component
 638 of its tuning curve given as

$$639 \quad f_\theta = \left\| \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{is} r_\theta(s) ds \right\| \#(24)$$

640 where $r_\theta(s)$ is the activity of the neuron that is selective to θ at the end of the delay period of a trial
 641 stimulated at s .

642 We calculated the mean and standard deviation across neurons.

$$643 \quad \text{mean}(f) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f_\theta d\theta \#(25)$$

$$644 \quad \text{std}(f) = \sqrt{\frac{1}{2\pi} \int_{-\pi}^{\pi} (f_\theta - \text{mean}(f))^2 d\theta} \quad (26)$$

645 The normalized std (std/mean) was used to quantify translation-invariance in Fig. 6-11.

646 Table of parameters for spatially structured network

| Parameter | Description | Value |
|----------------------------|--|-----------|
| N | Number of population in each group | 64 |
| τ_E | Time constant of excitatory neurons | 20 |
| τ_I | Time constant of inhibitory neurons | 10 |
| τ_{EE} | Time constant of E-to-E synapses | 100 |
| τ_{EI} | Time constant of I-to-E synapses | 10 |
| τ_{IE} | Time constant of E-to-I synapses | 25 |
| τ_{II} | Time constant of I-to-I synapses | 10 |
| τ_o | Time constant of external stimulus | 100 |
| J_{EE} | Amplitude of E-to-E synaptic weight | 100 |
| J_{EI} | Amplitude of I-to-E synaptic weight | 100 |
| J_{IE} | Amplitude of E-to-I synaptic weight | 200 |
| J_{II} | Amplitude of I-to-I synaptic weight | 200 |
| J_o | Amplitude of external stimulus | 270 |
| σ_{EE}, σ_{IE} | Width of excitatory synaptic connections | 0.2π |
| σ_{EI}, σ_{II} | Width of inhibitory synaptic connections | 0.1π |
| σ_o | Width of stimulus | 0.25π |
| h_0 | Baseline of stimulus | 200 |
| p | 1 - perturbation strength | vary |
| α_d | Learning rate of differential rule | 1e-5 |
| α_h | Learning rate of homeostatic rule | 2e-8 |
| r_0 | Target firing rate of homeostatic rule | 20 |
| t_{stim} | Stimulation duration | 500 |
| t_{total} | Stimulation plus delay period | 3500 |

647

648 Acknowledgments

649 The authors acknowledge the support from NYU-ECNU Institute of Brain and Cognitive Science

650 References

- 651 1. Knierim JJ, Zhang K. Attractor dynamics of spatially correlated neural activity in the limbic system.
652 Annu Rev Neurosci. 2012;35:267–85.
- 653 2. Goldman MS, Compte A, Wang XJ. Neural Integrator Models. In: Encyclopedia of Neuroscience.
654 Elsevier Ltd; 2009. p. 165–78.
- 655 3. Durstewitz D, Seamans JK, Sejnowski TJ. Neurocomputational Models of Working Memory. Nat
656 Neurosci. 2000;3(11s):1184–91.

- 657 4. Aksay E, Baker R, Seung HS, Tank DW. Correlated Discharge among Cell Pairs within the
658 Oculomotor Horizontal Velocity-to-Position Integrator. *J Neurosci*. 2003;
- 659 5. Yoon K, Buice MA, Barry C, Hayman R, Burgess N, Fiete IR. Specific evidence of low-dimensional
660 continuous attractor dynamics in grid cells. *Nat Neurosci*. 2013 Aug 14;16(8):1077–84.
- 661 6. Wimmer K, Nykamp DQ, Constantinidis C, Compte A. Bump attractor dynamics in prefrontal
662 cortex explains behavioral precision in spatial working memory. *Nat Neurosci*. 2014;17(3):431–9.
- 663 7. Brody CD, Romo R, Kepecs A. Basic mechanisms for graded persistent activity: discrete attractors,
664 continuous attractors, and dynamic representations. *Curr Opin Neurobiol*. 2003 Apr;13(2):204–
665 11.
- 666 8. Seung HS. How the brain keeps the eyes still. *Proc Natl Acad Sci U S A*. 1996;93(23):13339–44.
- 667 9. Koulakov AA, Raghavachari S, Kepecs A, Lisman JE. Model for a robust neural integrator. *Nat*
668 *Neurosci*. 2002;5(8):775–82.
- 669 10. Goldman MS, Levine JH, Major G, Tank DW, Seung HS. Robust Persistent Neural Activity in a
670 Model Integrator with Multiple Hysteretic Dendrites per Neuron. *Cereb Cortex*. 2003;
- 671 11. Lim S, Goldman MS. Balanced cortical microcircuitry for maintaining information in working
672 memory. *Nat Neurosci*. 2013 Sep 18;16(9):1306–14.
- 673 12. Lim S, Goldman MS. Balanced Cortical Microcircuitry for Spatial Working Memory Based on
674 Corrective Feedback Control. *J Neurosci*. 2014;34(20):6790–806.
- 675 13. Itskov V, Hansel D, Tsodyks M. Short-term facilitation may stabilize parametric working memory
676 trace. *Front Comput Neurosci*. 2011 Oct 24;5.
- 677 14. Seeholzer A, Deger M, Gerstner W. Stability of working memory in continuous attractor networks
678 under the control of shortterm plasticity. Burak Y, editor. *PLoS Comput Biol*. 2019 Apr
679 19;15(4):e1006928.
- 680 15. Arnold DB, Robinson DA. A neural network model of the vestibulo-ocular reflex using a local
681 synaptic learning rule. *Philos Trans R Soc Lond B Biol Sci*. 1992;337(1281):327–30.
- 682 16. Major G, Baker R, Aksay E, Mensh B, Seung HS, Tank DW. Plasticity and tuning by visual feedback
683 of the stability of a neural integrator. *Proc Natl Acad Sci U S A*. 2004;101(20):7739–44.
- 684 17. MacNeil D, Eliasmith C. Fine-tuning and the stability of recurrent neural networks. Vasilaki E,
685 editor. *PLoS One*. 2011 Sep 27;6(9):e22885.
- 686 18. Xie X, Seung HS. Spike-based learning rules and stabilization of persistent neural activity. In: Solla
687 SA, Leen TK, Müller K, editors. *Advances in Neural Information Processing Systems*. MIT Press;
688 2000. p. 199–205.
- 689 19. Nygren E, Ramirez A, McMahan B, Aksay E, Senn W. Learning temporal integration from internal
690 feedback. *bioRxiv*. 2019;
- 691 20. Federer C, Zylberberg J. A self-organizing short-term dynamical memory network. *Neural*
692 *Networks*. 2018 Oct 1;106:30–41.
- 693 21. Kosko B. Differential Hebbian learning. In: *AIP Conference Proceedings*. AIP; 1986. p. 277–82.

- 694 22. Der R, Martius G. Novel plasticity rule can explain the development of sensorimotor intelligence.
695 Proc Natl Acad Sci. 2015 Nov 10;112(45):E6224–32.
- 696 23. Roberts PD. Computational consequences of temporally asymmetric learning rules: I. Differential
697 Hebbian learning. J Comput Neurosci. 1999;
- 698 24. Harry Klopf A. A neuronal model of classical conditioning. Psychobiology. 1988;
- 699 25. Wörgötter F, Porr B. Temporal Sequence Learning, Prediction, and Control: A Review of Different
700 Models and Their Relation to Biological Mechanisms. Neural Comput. 2005;17(2):245–319.
- 701 26. Gluck MA, Parker DB, Reifsnider E. Erratum to: Some biological implications of a differential-
702 Hebbian learning rule. Vol. 17, Psychobiology. 1989. p. 110–110.
- 703 27. Turrigiano GG, Leslie KR, Desai NS, Rutherford LC, Nelson SB. Activity-dependent scaling of
704 quantal amplitude in neocortical neurons. Nature. 1998;391(6670):892–6.
- 705 28. Van Rossum MCW, Bi GQ, Turrigiano GG. Stable Hebbian learning from spike timing-dependent
706 plasticity. J Neurosci. 2000 Dec 1;20(23):8812–21.
- 707 29. Renart A, Song P, Wang X-J. Robust spatial working memory through homeostatic synaptic
708 scaling in heterogeneous cortical networks. Neuron. 2003;38(3):473–85.
- 709 30. Pool RR, Mato G. Hebbian Plasticity and Homeostasis in a Model of Hypercolumn of the Visual
710 Cortex. Neural Comput. 2010 Jul 27;1859(7):1837–59.
- 711 31. Romo R, Brody CD, Hernández A, Lemus L. Neuronal correlates of parametric working memory in
712 the prefrontal cortex. Nature. 1999 Jun;399(6735):470–3.
- 713 32. Machens CK, Romo R, Brody CD. Flexible control of mutual inhibition: A neural model of two-
714 interval discrimination. Science (80-). 2005 Feb 18;307(5712):1121–4.
- 715 33. Wu S, Wong KYM, Fung CCA, Mi Y, Zhang W. Continuous Attractor Neural Networks: Candidate of
716 a Canonical Model for Neural Information Representation. F1000Research. 2016 Feb 10;5:156.
- 717 34. Wang M, Yang Y, Wang C-J, Gamo NJ, Jin LE, Mazer JA, et al. NMDA Receptors Subserve
718 Persistent Neuronal Firing during Working Memory in Dorsolateral Prefrontal Cortex. Neuron.
719 2013 Feb;77(4):736–49.
- 720 35. Zenke F, Agnes EJ, Gerstner W. Diverse synaptic plasticity mechanisms orchestrated to form and
721 retrieve memories in spiking neural networks. Nat Commun. 2015 Nov 21;6(1):6922.
- 722 36. Litwin-Kumar A, Doiron B. Formation and maintenance of neuronal assemblies through synaptic
723 plasticity. Nat Commun. 2014 Dec 14;5(1):5319.
- 724 37. Goldman-Rakic P. Cellular basis of working memory. Neuron. 1995 Mar 1;14(3):477–85.
- 725 38. Rao SG, Williams G V., Goldman-Rakic PS. Isodirectional Tuning of Adjacent Interneurons and
726 Pyramidal Cells During Working Memory: Evidence for Microcolumnar Organization in PFC. J
727 Neurophysiol. 1999 Apr 1;81(4):1903–16.
- 728 39. Constantinidis C, Goldman-Rakic PS. Correlated Discharges Among Putative Pyramidal Neurons
729 and Interneurons in the Primate Prefrontal Cortex. J Neurophysiol. 2002 Dec 1;88(6):3487–97.
- 730 40. Strang G. Introduction to Linear Algebra. Wellesley-Cambridge Press; 2017.

- 731 41. Rotaru DC, Yoshino H, Lewis DA, Ermentrout GB, Gonzalez-Burgos G. Glutamate Receptor
732 Subtypes Mediating Synaptic Activation of Prefrontal Cortex Neurons: Relevance for
733 Schizophrenia. *J Neurosci*. 2011 Jan 5;31(1):142–56.
- 734 42. Constantinidis C, Franowicz MN, Goldman-Rakic PS. Coding specificity in cortical microcircuits: A
735 multiple-electrode analysis of primate prefrontal cortex. *J Neurosci*. 2001;21(10):3646–55.
- 736 43. Akil AE, Rosenbaum R, Josić K. Synaptic Plasticity in Correlated Balanced Networks. *bioRxiv*. 2020
737 Apr 26;
- 738 44. Legenstein R, Maass W. Branch-Specific Plasticity Enables Self-Organization of Nonlinear
739 Computation in Single Neurons. *J Neurosci*. 2011 Jul 27;31(30):10787–802.
- 740 45. Ocker GK, Litwin-Kumar A, Doiron B. Self-Organization of Microcircuits in Networks of Spiking
741 Neurons with Plastic Synapses. Latham PE, editor. *PLOS Comput Biol*. 2015 Aug
742 20;11(8):e1004458.
- 743 46. Mongillo G, Curti E, Romani S, Amit DJ. Learning in realistic networks of spiking neurons and
744 spike-driven plastic synapses. *Eur J Neurosci*. 2005 Jun;21(11):3143–60.
- 745 47. Zenke F, Gerstner W, Ganguli S. The temporal paradox of Hebbian learning and homeostatic
746 plasticity. Vol. 43, *Current Opinion in Neurobiology*. Elsevier Ltd; 2017. p. 166–76.
- 747 48. Chen X, Bialek W. Searching for long time scales without fine tuning. *arxiv*. 2020 Aug 26;1–13.
- 748 49. Machens CK, Romo R, Brody CD. Functional, But Not Anatomical, Separation of “What” and
749 “When” in Prefrontal Cortex. *J Neurosci*. 2010 Jan 6;30(1):350–60.
- 750 50. Murray JD, Bernacchia A, Roy NA, Constantinidis C, Romo R, Wang XJ. Stable population coding
751 for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proc Natl
752 Acad Sci U S A*. 2017 Jan 10;114(2):394–9.
- 753 51. Inagaki HK, Inagaki M, Romani S, Svoboda K. Low-Dimensional and Monotonic Preparatory
754 Activity in Mouse Anterior Lateral Motor Cortex. *J Neurosci*. 2018 Apr 25;38(17):4163–85.
- 755 52. Cueva CJ, Saez A, Marcos E, Genovesio A, Jazayeri M, Romo R, et al. Low-dimensional dynamics
756 for working memory and time encoding. *Proc Natl Acad Sci U S A*. 2020 Sep 15;117(37):23021–
757 32.
- 758 53. Druckmann S, Chklovskii DB. Neuronal circuits underlying persistent representations despite time
759 varying activity. *Curr Biol*. 2012;22(22):2095–103.
- 760 54. Alemi A, Denève S, Machens CK, Slotine JJ. Learning nonlinear dynamics in efficient, balanced
761 spiking networks using local plasticity rules. In: *AAAI Conference*. 2018. p. 588–95.
- 762 55. Savin C, Triesch J. Emergence of task-dependent representations in working memory circuits.
763 *Front Comput Neurosci*. 2014;8(MAY):1–12.
- 764 56. Hertz J, Krogh A, Palmer RG, Horner H. Introduction to the theory of neural computation. *Phys
765 Today*. 1991;44(12):70.
- 766 57. Chen Z, Haykin S, Eggermont JJ, Becker S. Correlative learning: a basis for brain and adaptive
767 systems. Vol. 49. John Wiley & Sons; 2008.
- 768 58. Lim S. Hebbian learning revisited and its inference underlying cognitive function. *Curr Opin Behav*

- 769 Sci. 2021;38:96–102.
- 770 59. Vogels TP, Sprekeler H, Zenke F, Clopath C, Gerstner W. Inhibitory Plasticity Balances Excitation
771 and Inhibition in Sensory Pathways and Memory Networks. *Science* (80-). 2011 Dec
772 16;334(6062):1569–73.
- 773 60. Froemke RC. Plasticity of Cortical Excitatory-Inhibitory Balance. *Annu Rev Neurosci*.
774 2015;38(1):195–219.
- 775 61. Luz Y, Shamir M. Balancing feed-forward excitation and inhibition via hebbian inhibitory synaptic
776 plasticity. *PLoS Comput Biol*. 2012 Jan;8(1).
- 777

778 Supplementary Information

779 **Figure S1, related to Fig. 6, Elongation of time constant associated with each eigenvector similar to** 780 **Fourier modes under differential plasticity.**

781 A: Time scale of each Fourier mode. For each Fourier mode, a time constant was estimated by projecting
782 population activity onto a sinusoid of different frequencies and fitting the time course with exponential
783 decay. The negative reciprocals of these time constants have good correspondence with the eigenvalues
784 shown in Fig 6D. B: Eigenvectors related to eigenvalues in Fig. 6D during the evolution of learning
785 dynamics. The real part of the eigenvectors corresponding to the first, third, and fifth leading
786 eigenvalues (even ones omitted because of redundancy) is plotted. The shape of the eigenvectors is
787 close to sinusoids, suggesting preservation of translation-invariance.

788

789 **Figure S2, related to Fig. 11. The effect of combined plasticity under larger local perturbation.**

790 A-B: same as Fig 11 B-C, but under 30% local perturbation, as indicated by arrowhead in C-D. C-D: copy
791 of Fig 11E-F. For larger local perturbation, translation-invariance breaks down while decoding errors are
792 low. The distortion of activity patterns is dissimilar to that only with differential plasticity because
793 homeostatic plasticity keeps neurons from falling silent in some trials but not in others.

794

795 **Figure S3. The effect of differential, homeostatic, and combined plasticity rules in positive feedback** 796 **network under 30% pre-synaptic perturbation.**

797 A-C: Decoding errors for eight stimulus groups during the recovery. Under the local perturbation,
798 translation-invariance breaks down in the positive feedback models, and the bump activity tends to drift
799 as in a negative derivative feedback network. The effects under pre-synaptic perturbations were only
800 shown for simplicity. D: Activity pattern with homeostatic plasticity at around trial 4000 (arrows in A). E,
801 F: Activity pattern with differential and combined plasticity at around trial 2000 (arrows in B, C). As in
802 the negative derivative feedback network, homeostatic synaptic plasticity is not effective under pre-
803 synaptic perturbation (D). In contrast, differential plasticity effectively stops the drift (E), as well as the
804 combined plasticity (F). Note that the activity patterns under the differential (B) and combined plasticity
805 (C) start to break down around 3000 trials after the recovery. This breakdown may arise from

806 implementing a raw dr/dt in differential plasticity, which is sensitive to a slight change of activity pattern
807 as it evolves from stimulus-evoked patterns to the stereotypical delay activity patterns.

808

Figure 2

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.17.444447>; this version posted May 17, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

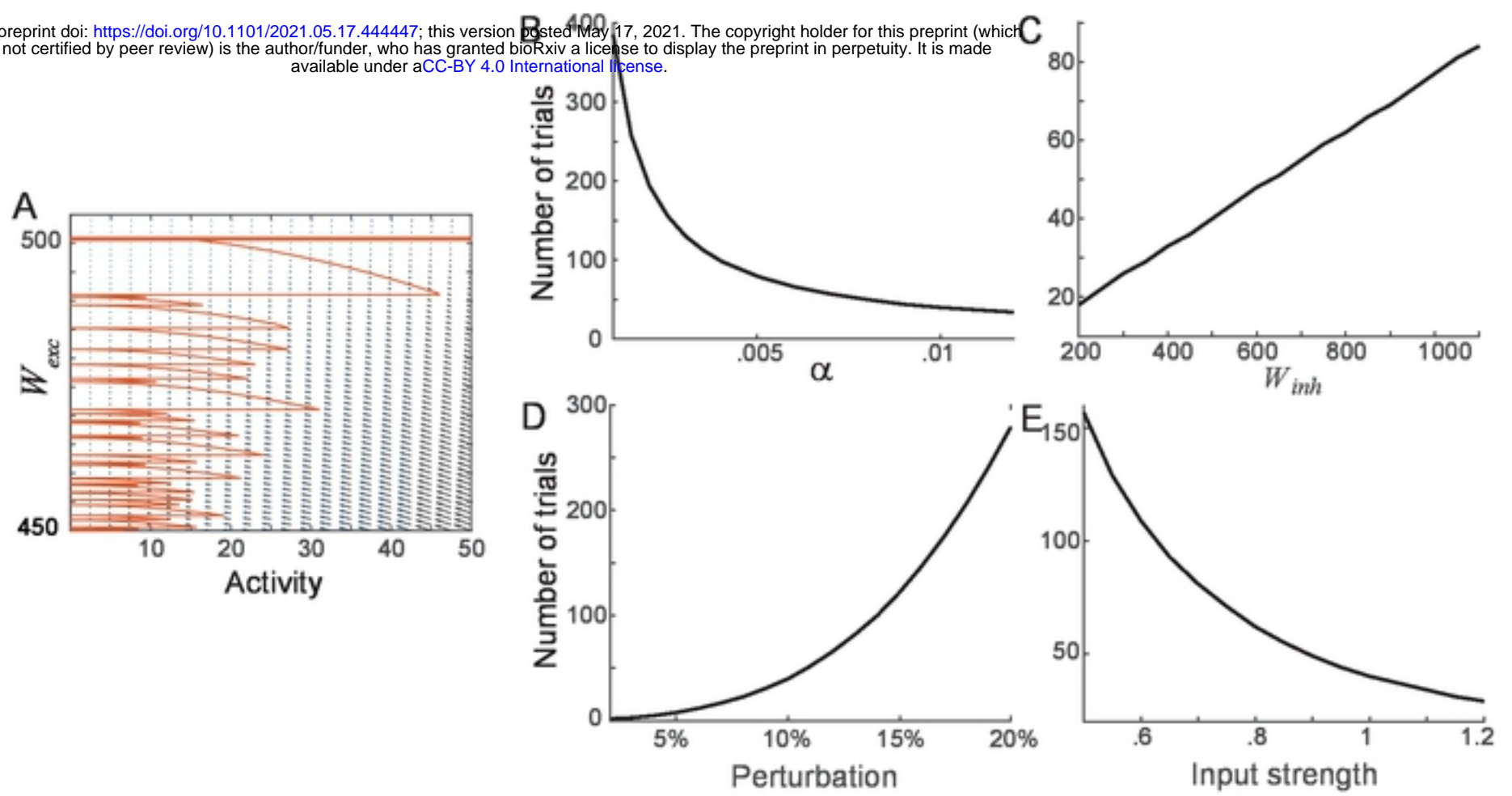


Figure 3

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.17.444447>; this version posted May 17, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

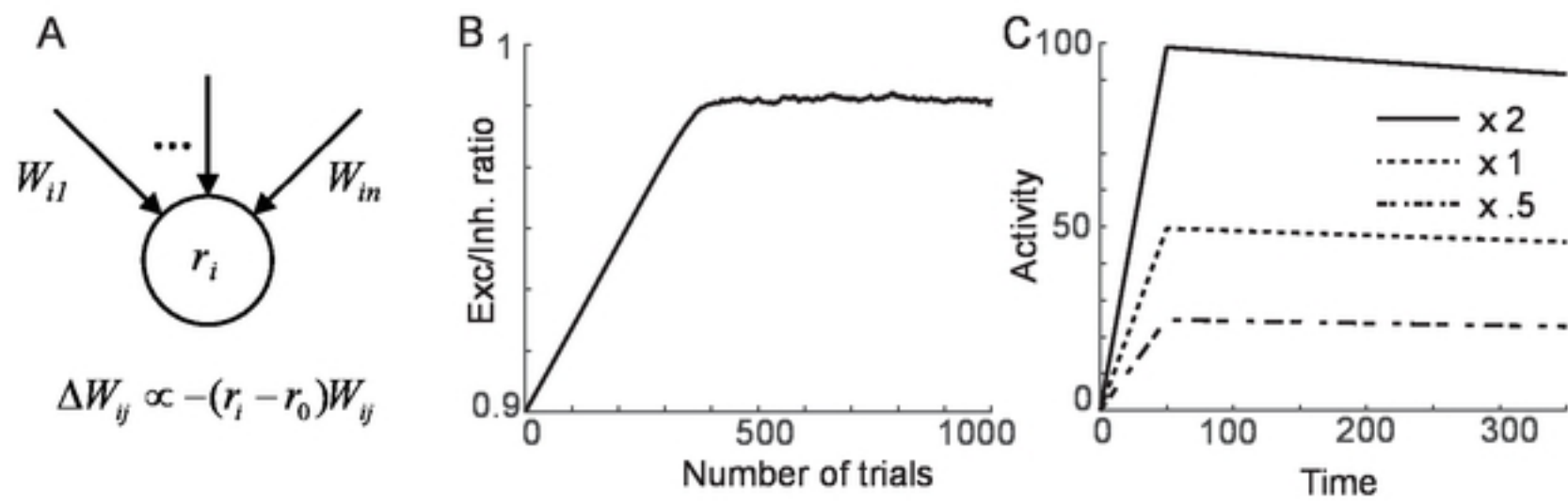


Figure 4

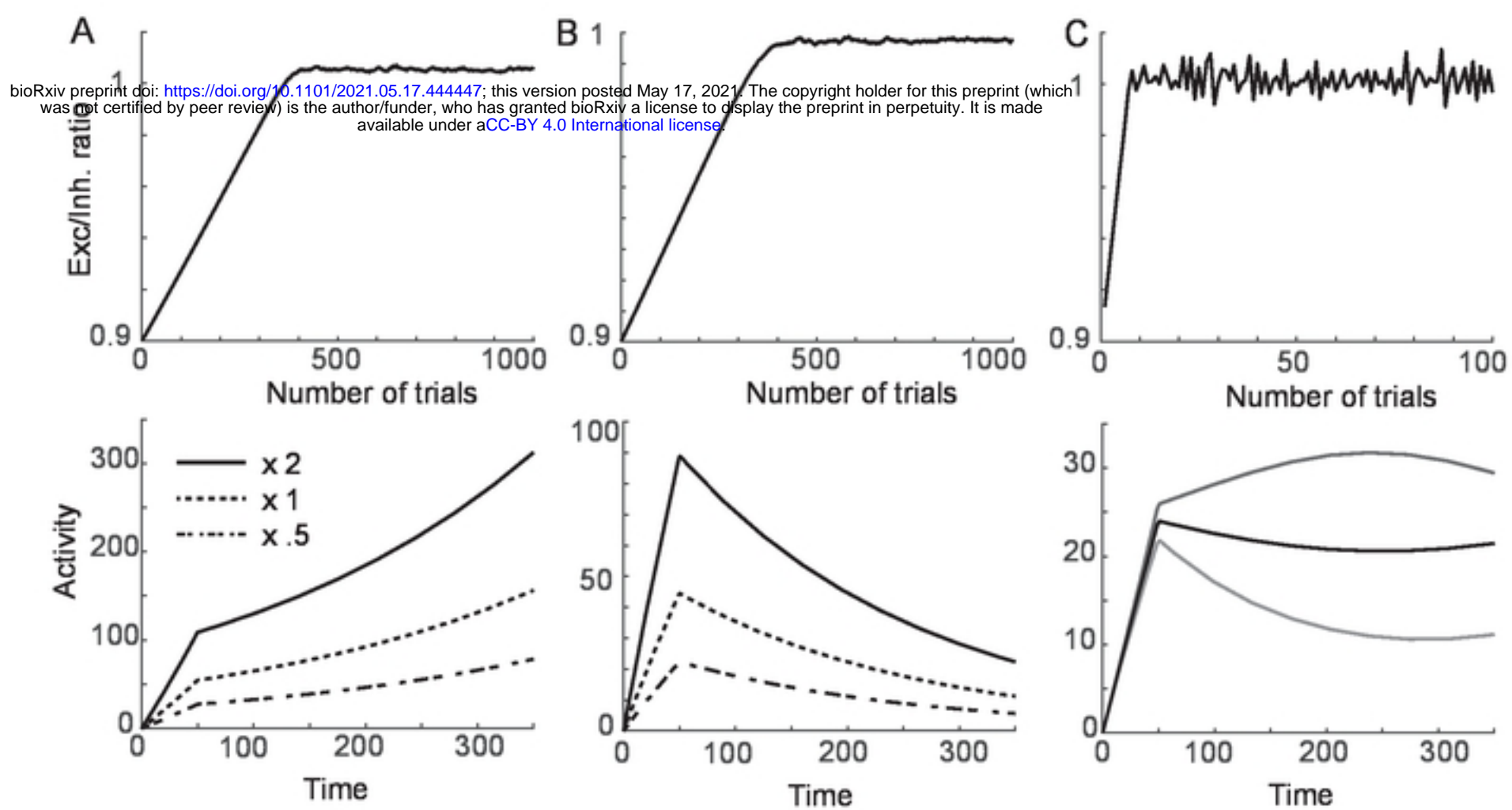


Figure 5

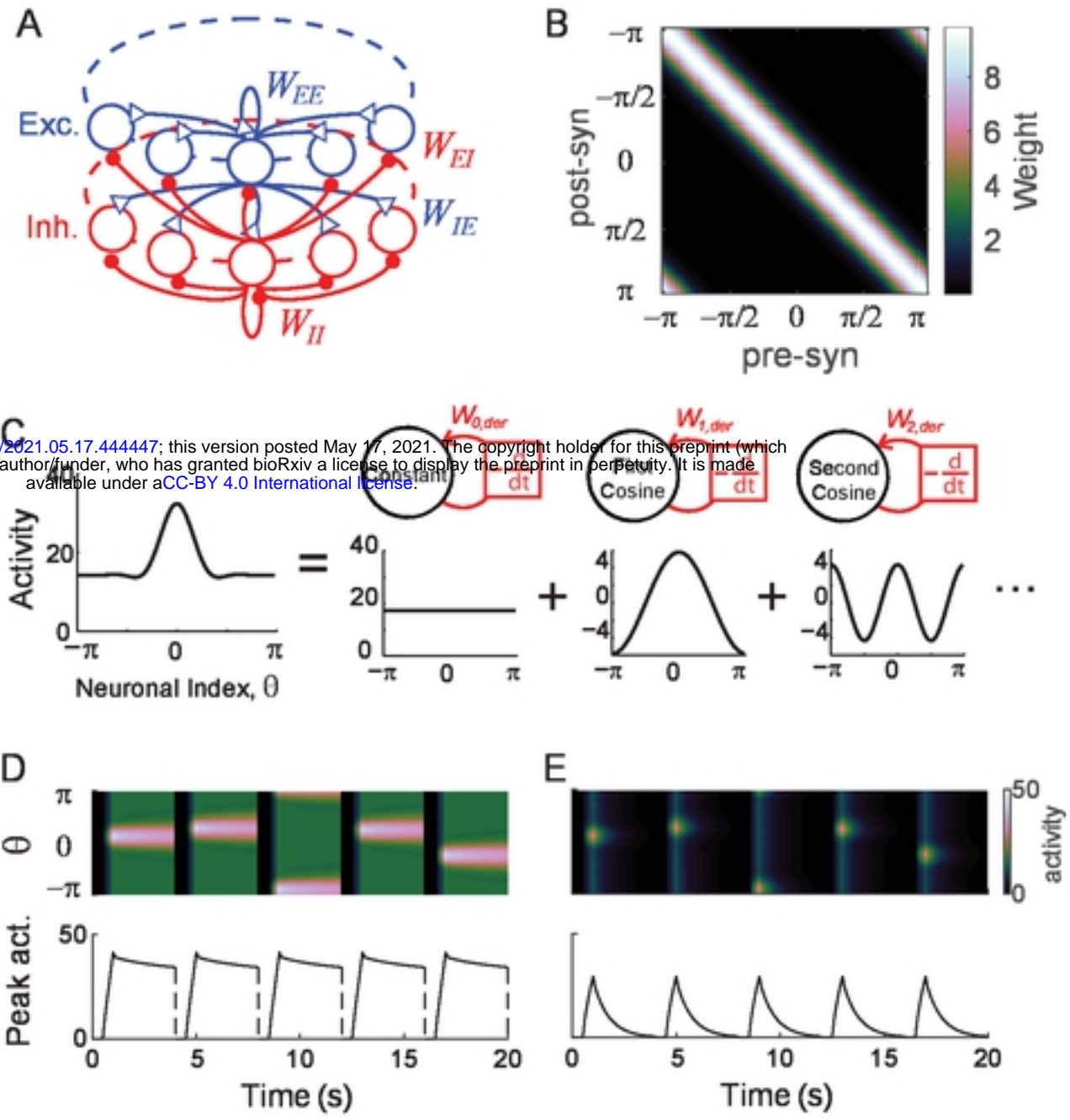


Figure 6

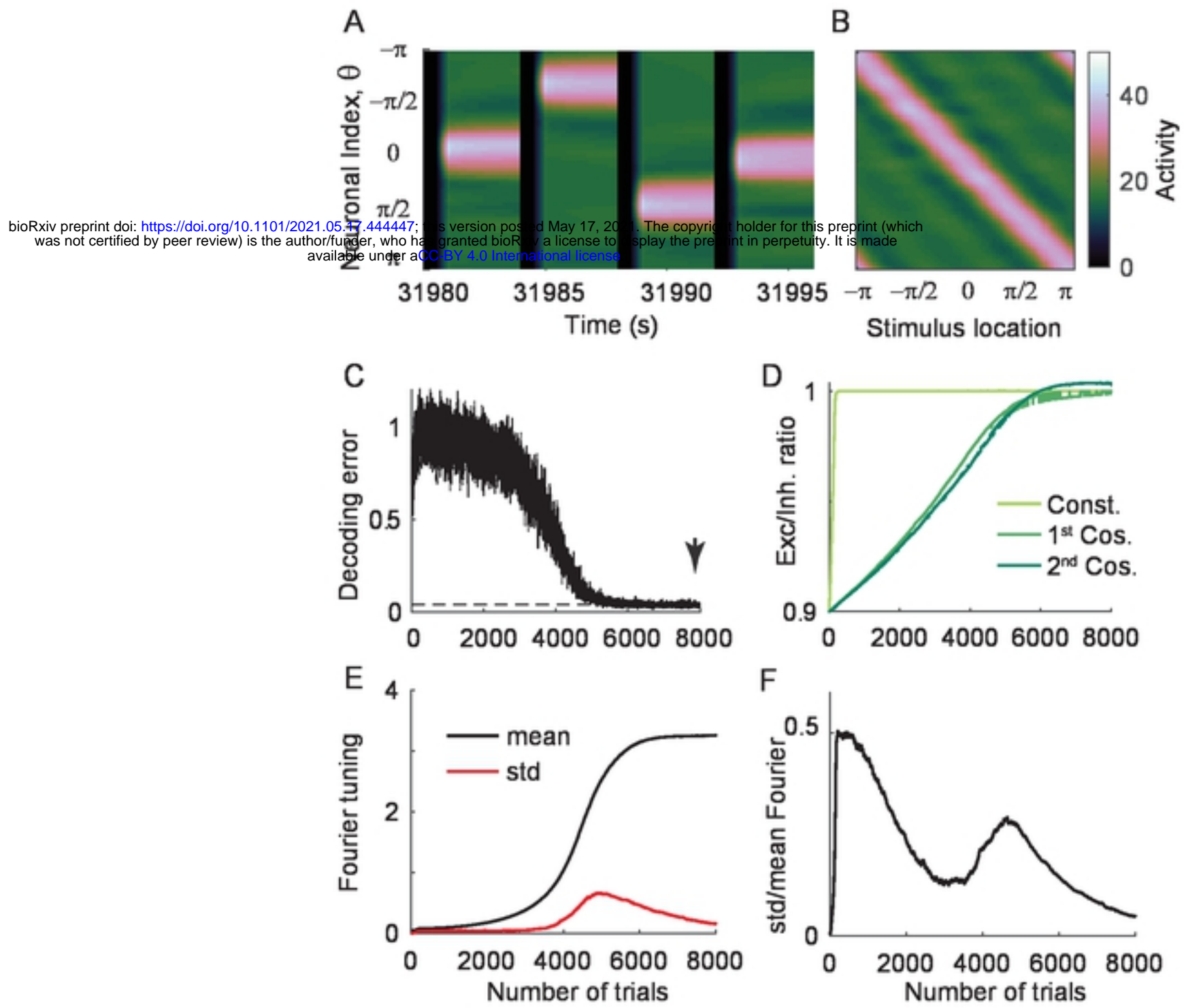


Figure 7

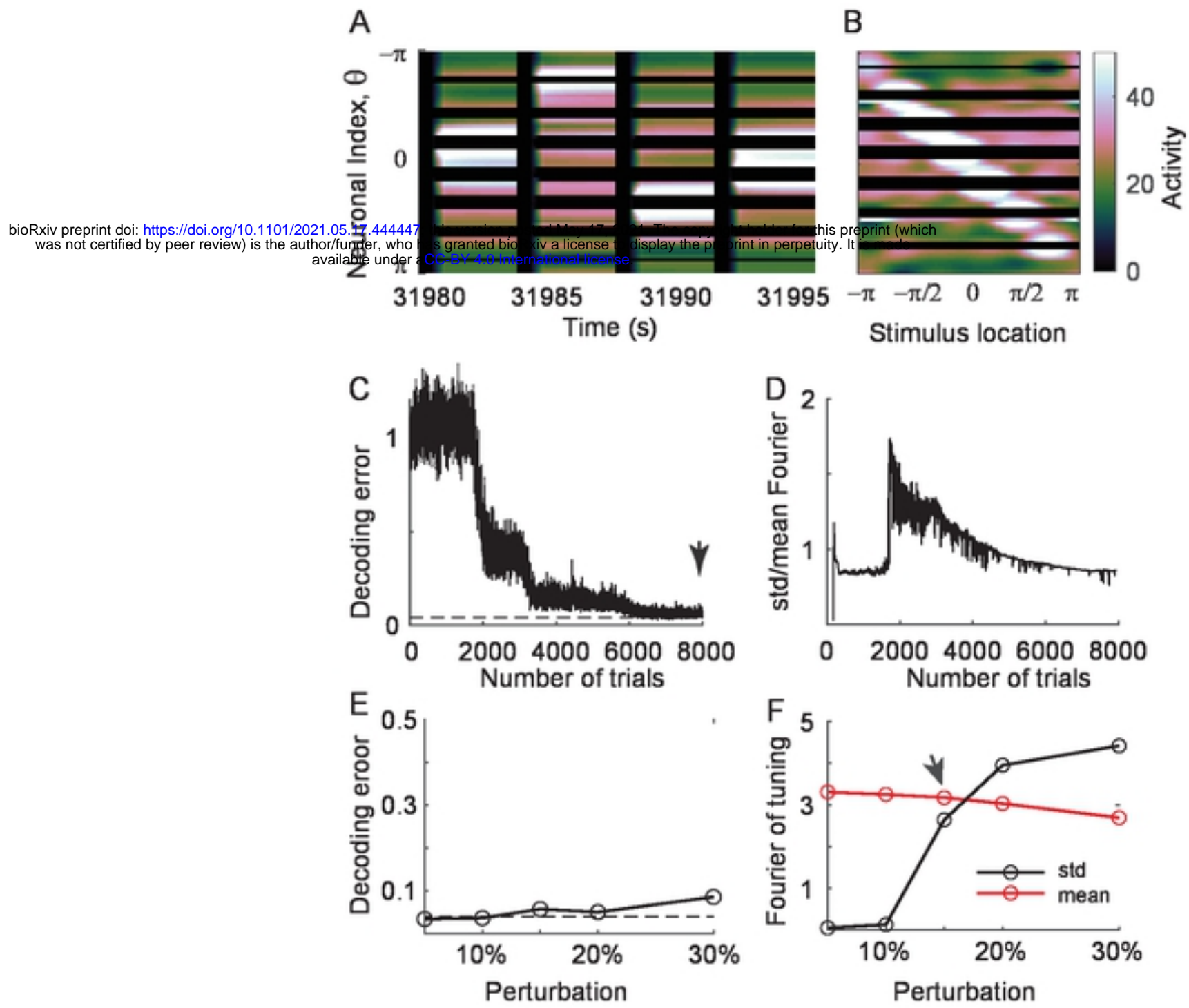


Figure 8

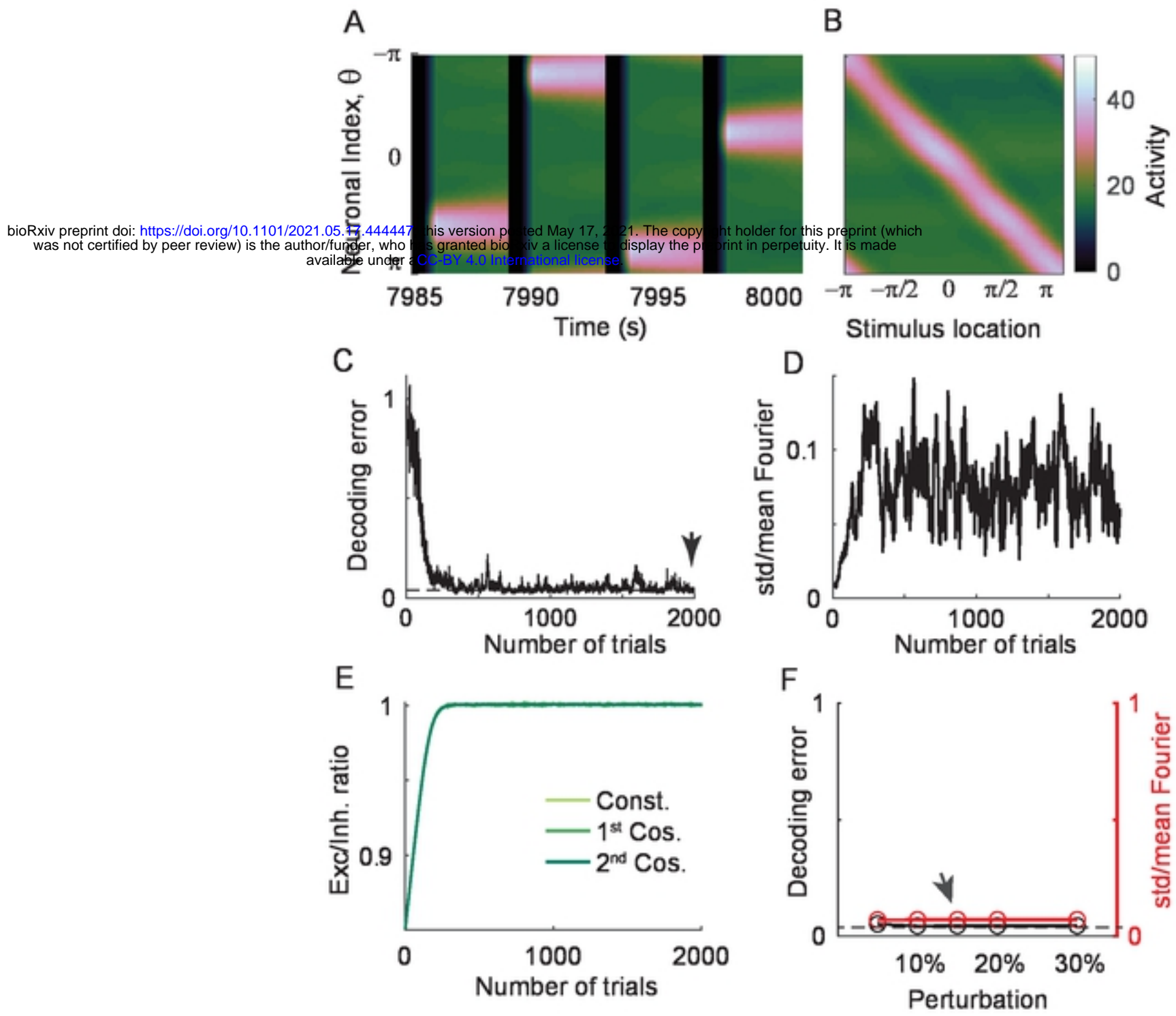


Figure 9

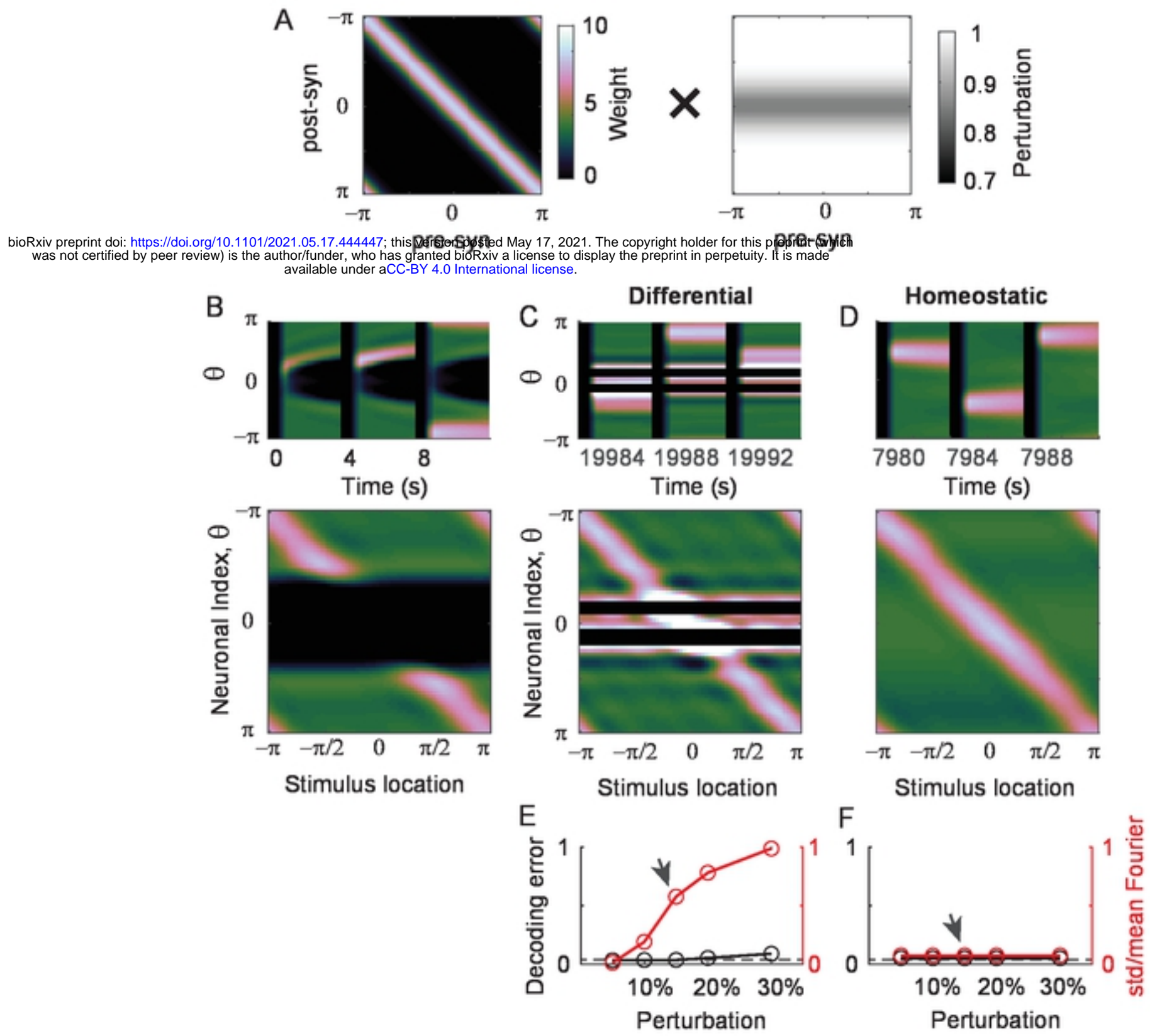


Figure 10

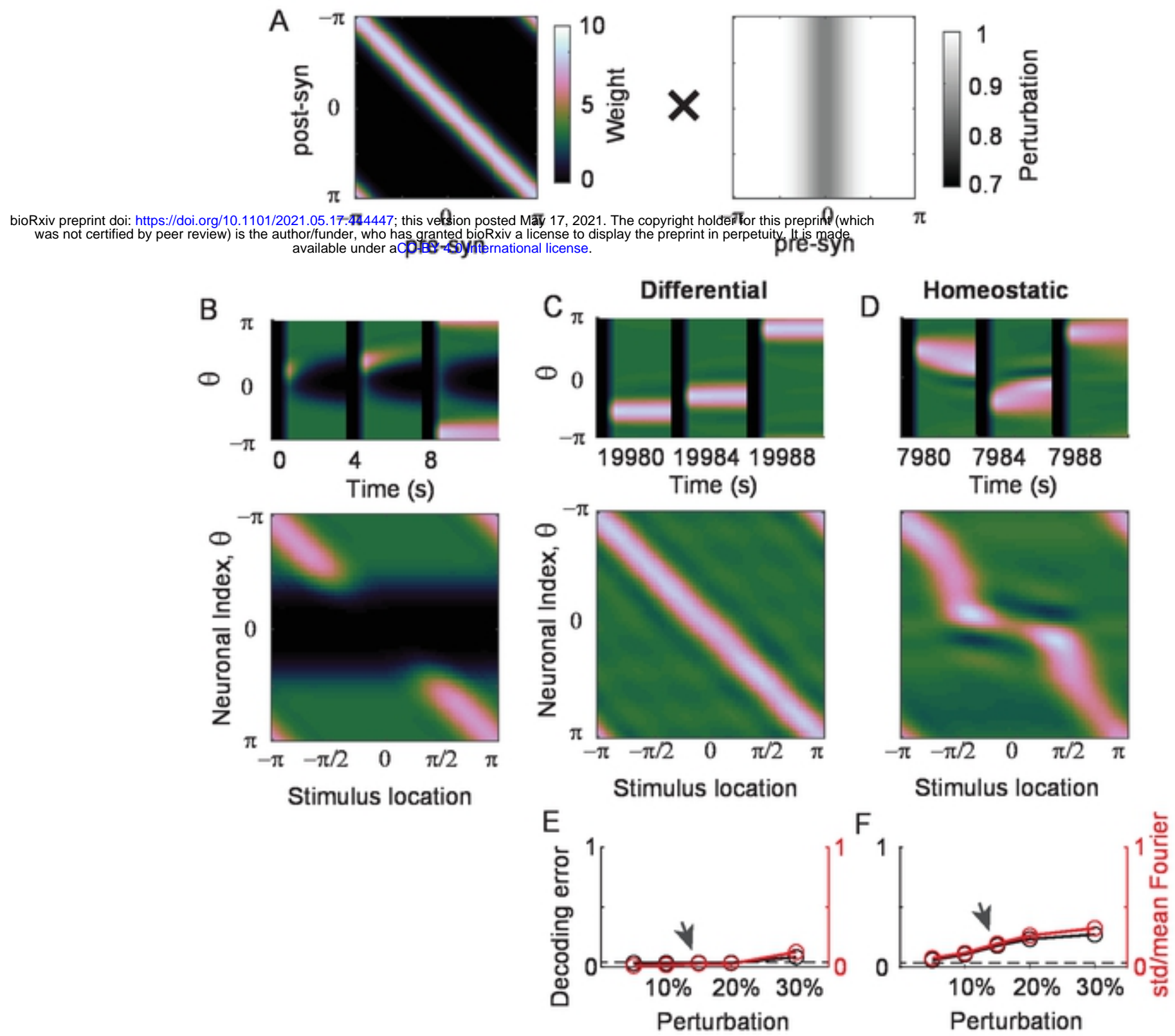


Figure 11

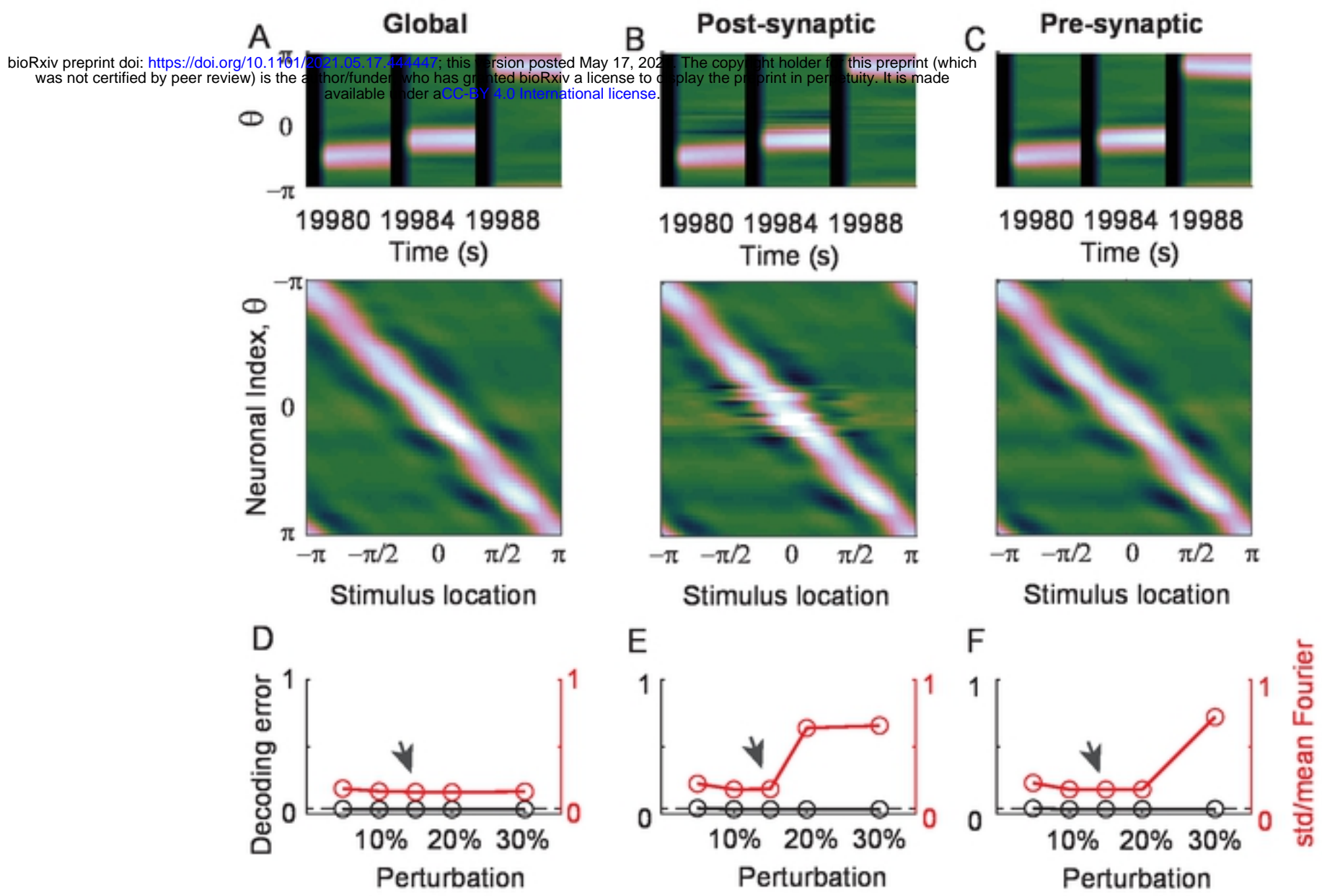


Figure 1

