
1

2 **Modification of a Reward-Modulated Hebbian**

3 **Learning Rule as a Model of Working Memory**

4 **Emergence**

5

6

7 **Tammy Tran Brad Theilman Nuttida Rungratsameetaweemana**

8 Neurosciences Graduate Program, University of California, San Diego

9 ttt075@ucsd.edu, btheilma@ucsd.edu, nrungrat@ucsd.edu

10

11

12 **Abstract**

13 In liquid state machines with generic cortical microcircuits, synaptic plasticity

14 can be optimized by reward-modulated Hebbian learning, eliminating the need

15 for supervised learning (Hoerzer, Legenstein, and Maass, 2014). Reward-

16 modulated Hebbian learning can thus lead to autonomous emergence of task-

17 specific working memory during the learning of computational rules. However,

18 in Hoerzer, Legenstein, and Maass (2014), reward-modulated Hebbian learning

19 was modeled with an all-or-none modulatory signal that permitted synaptic

20 weight change only above a criterion level of learning. We use liquid computing

21 models to investigate working memory emergence via Hebbian learning with

22 non-binary modulatory signals. We implement a nonbinary, but discrete

23 modulatory signal and an analog signal. In doing so, we model physiological

24 conditions of tonic and phasic output of reward-mediating systems like the

25 dopaminergic system. We find that the effects of analog modulatory signals on

26 working memory emergence improve reward-modulated Hebbian learning in

27 liquid state machines. We propose that reward-modulated Hebbian learning in

28 generic microcircuits of neurons can abstractly model general cognitive

29 processes.

30

31 **1 Introduction**

32 The liquid state machine (LSM) is a conceptual framework that assumes generic recurrent neural

33 microcircuits where neurons are randomly connected to one another. Specifically, the vector of

34 contributions of all the neurons in the microcircuit to the membrane potential at time t of a readout

35 neuron is referred to as the liquid state $x(t)$, and this is all the information about the circuit state a

36 readout neuron has access to [3]. The LSMs do not require task-dependent constructions and need

37 not be engineered for a specific task, and hence this framework can be used to investigate a wide

38 range of computational tasks. The liquid state of an LSM is assumed to vary continuously over

39 time and to be sufficiently sensitive to information needed for specific tasks. In addition to this

40 universal computing power, the LSM framework has a capability of turning time-varying circuit

41 inputs into spatio-temporal activity pattern that represents the circuit dynamics. These

42 characteristics make it possible for researchers to train LSMs to fulfill a large variety of complex

43 computational tasks.

44 While past studies have shown that the LSMs can successfully be trained to learn several different
45 tasks, the training paradigms generally provide the microcircuits with knowledge of the desired
46 activity output (supervised learning). This type of learning also presupposes another neural
47 network that is capable of a particular computational task used, and thus cannot explain how
48 specialization first emerges. Additionally, the feedback provided in supervised learning represents
49 global activity of the entire network rather than localized dynamical activity of specific neurons,
50 resulting in physiologically inaccurate network outputs. Due to these pitfalls of supervised
51 learning, the present study examines working memory emergence by training recurrent neural
52 networks through unsupervised, reward-modulated Hebbian learning, where feedback provided to
53 the neural network represents local activity between the pre- and postsynaptic neurons.

54

55 2 Methods

56 To test working memory emergence via reward-modulated Hebbian learning in recurrent neural
57 networks, a generic network model was implemented as previously described [1]. Leaky integrator
58 neurons ($N = 1000$) were sparsely, recurrently connected. Two external input streams $u_i(t)$ were
59 provided to the recurrent neurons, and the recurrent neurons provided output to a single readout
60 neuron, which in turn provided feedback to the recurrent network. All neurons were connected or
61 received input or feedback with probability p of 0.1. Synaptic weights between recurrent neurons
62 (W^{rec}) were randomly drawn from Gaussian distributions with zero mean and $1/(pN)$ variance.
63 Input and feedback weights (W^{in} and W^{fb}) were drawn from uniform distributions in $[-1, 1]$.
64 Output weights w to the readout neuron were initialized to zero and adjusted during training.

65 Membrane potentials $x_j(t)$ per recurrent neuron j were initialized to zero, and network dynamics
66 per simulation time-step $dt = 1ms$ is given by

$$67 \quad \frac{\tau dx_i}{dt} = -x_i(t) + \lambda \sum_{i=1}^N W_{ij}^{rec} r_j(t) + \sum_{i=1}^M W_{ij}^{in} u_j(t) + \sum_{i=1}^L W_{ij}^{fb} z_j(t)$$

68 per recurrent neuron i . The chaoticity level λ was 1.2, and the network constant τ was 50ms.

69 The firing rate $r_i(t)$ per recurrent neuron j is given by $r_j(t) = \tanh[x_j(t)] + \zeta_j^{state}(t)$, with the
70 noise $\zeta_j^{state}(t)$ drawn from uniform distributions in $[-0.05, 0.05]$.

71 Output at the readout neuron is given by $z(t) = \mathbf{w}^T \mathbf{r}(t) + \zeta(t)$, where $\mathbf{r}(t)$ is the column vector of
72 firing rates $r_j(t)$ and the zero-mean exploration noise $\zeta(t)$ was drawn independently at each time
73 step from uniform distributions in $[-0.5, 0.5]$.

74 Weight change for $w(t)$ for each dt is given by $\Delta w(t) = \eta(t)[z(t) - z_{avg}(t)]M(t)r(t)$.

75 $\eta(t) = \frac{\eta_{init}}{1+\frac{t}{T}}$ is a linearly decaying learning rate with $\eta_{init} = 0.0005$ and $T = 20s$. $z_{avg}(t)$ is the

76 average readout output given by $z_{avg}(t) = \left(1 - \frac{dt}{\tau_{avg}}\right) z_{avg}(t - dt) + \left(\frac{dt}{\tau_{avg}}\right) z_i(t)$, with $\tau_{avg} =$
77 $5ms$. Initially, $M(t)$ is a binary modulatory signal of 1 if performance $P(t) > P_{avg}(t)$, 0 otherwise.
78 Performance $P(t)$ is given by $P(t) = -\sum_{i=1}^L [z(t) - f(t)]^2$ where $f(t)$ is the target output of the
79 readout neuron. $P_{avg}(t)$ is given by $P_{avg}(t) = \left(1 - \frac{dt}{\tau_{avg}}\right) P_{avg}(t - dt) + \left(\frac{dt}{\tau_{avg}}\right) P(t)$.

80 For the working memory computational task described previously, input streams $\widehat{u}_{on}(t)$ and
81 $\widehat{u}_{off}(t)$ were independently set to 1 with probability 0.0005, zero otherwise [1]. The pulses were
82 then smoothed to amplitude 0.4 and duration 100ms. Thus, $u_{on}(t) = \frac{1}{\sigma_u} (\theta_0 \circ (\widehat{u}_{on} * h)) * g$,

83 and $u_{off}(t)$ is similarly constructed with $\widehat{u}_{on}(t)$. The functions g and h are $g(s) = \exp(-s/\tau_L)\theta_1(s)$ and
84 $h(s) = \theta_1(s) - \theta_1(s - 100ms)$, with $\tau_L = 50ms$ and $\theta_x(s)$ being the Heaviside function equal to 0 for x
85 < 0 , x for $s = 0$, and 1 otherwise. The target function $f(t) = \frac{1}{\sigma_f} \widehat{f}(t) * g$ with amplitude 0.5 for the

86 readout neuron was then a smoothed version of $\widehat{f}(t) = 0.5$ if $\widehat{u}_{on}(t) = 1$, 0.5 if $\widehat{u}_{on}(t) = -1$, and
87 $\widehat{f}_i(t - dt)$ otherwise.

88 **3 Results**

89 **3.1 Reproduction of Hoezer, Legenstein, and Maass’s model**

90 This study began by recreating some of the main results from [1], namely, the production of
91 periodic signals and a simple working memory task. We created a 1000 neuron recurrent neural
92 network with dynamics identical to the model in [1] and trained the network to recreate a periodic
93 function consisting of the sum of five sine waves of differing frequencies. The network was
94 trained using the binary reward modulation signal for 500 seconds of simulation time. The figure
95 below compares the output of the network after training and the target function. It can be seen that
96 the network reproduces the target function closely, but over time the phase drifts away from the
97 perfectly periodic target function.

98

99

100

101

102

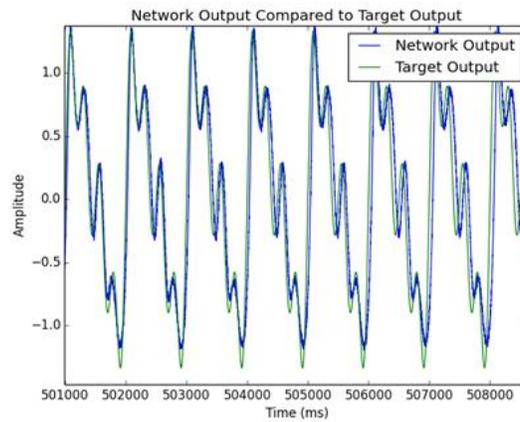
103

104

105

106

107



108 Figure 1. Comparison between the network output after training and the periodic target function
109 when using a binary modulatory signal

110

111 Next, we trained the same network to perform the simple working memory task from the [1]. Two
112 input streams were connected to the network with random, uniformly distributed weights to each
113 neuron in the network. The readout output was trained to drive to a high value (0.5) when the most
114 recent input stream that was active was input stream “on”, and to drive to a low output (-0.5) when
115 the most recent active input stream was stream “off”. The figure below shows the output of the
116 network compared to a target function that idealizes the desired outcome of the task:

117

118

119

120

121

122

123

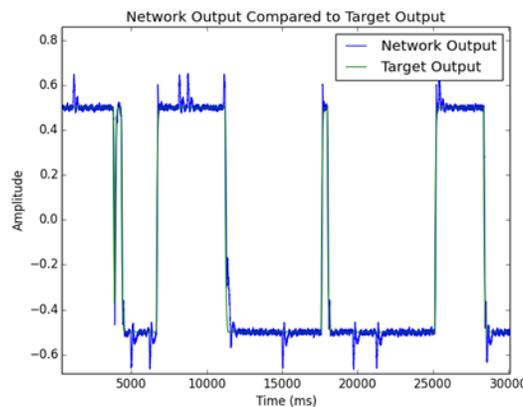
124

125

126

127

128



129 Figure 2. Comparison between the network output after training and the working memory target
130 function when using a binary modulatory signal

131 While the state transitions are clear and correct, the network transiently deviates from the target
132 function. These deviations are artifacts of the input stream activations and are present because
133 network dynamics are not instantaneous. Taken together, these results show that our model
134 successfully reproduces the key results of [1]. Thus, our model can be used as a platform upon
135 which to investigate how dopamine-like reward signals modulating Hebbian learning affect
136 learning in liquid state machines.

137

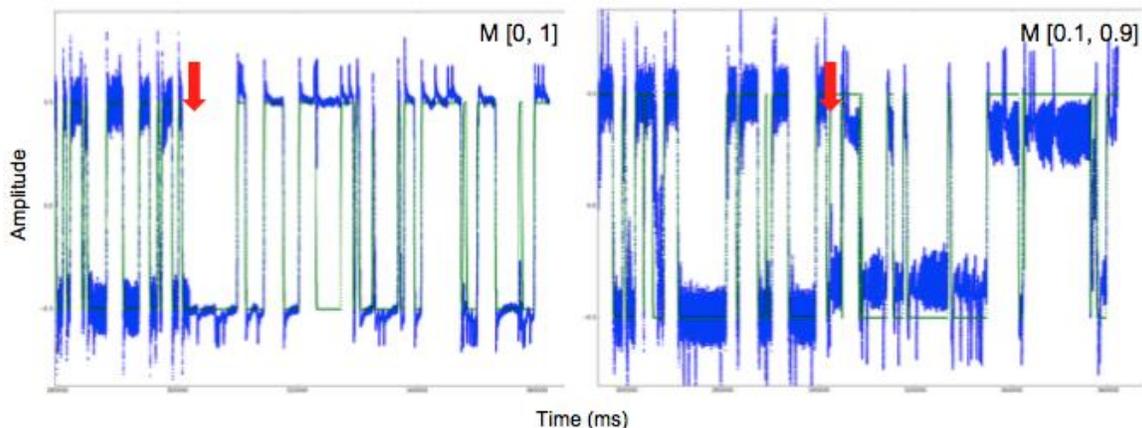
138 3.2 Modification of the modulatory signal

139 3.2.1 Implementing a nonbinary, but discrete modulatory signal

140 While $M(t) = 0$ or 1 distinguishes between absent and high transmission of a reward signal, a
141 binary signal does not accurately reflect physiological transmission of reward signals like
142 dopamine. In reality, dopamine transmission is both phasic and tonic, with some persistent,
143 baseline dopamine tone at all times. In addition, phasic dopamine transmission can vary in
144 amplitude, exhibiting characteristics of an analog signal. Therefore, to better represent
145 physiological reward signals, we altered $M(t)$ and tested the network's performance in the simple
146 working memory task from [1].

147 First, $M(t)$ was altered such that $M(t) = 0.9$ at performance greater than P_{avg} and 0.1 otherwise.
148 This nonbinary, yet discrete signal accurately reflects baseline reward signal or tonic dopamine
149 transmission, though without the continuous nature of an analog reward signal. With $M(t)$ defined
150 so, performance on the simple working memory task after 300s of training was measured.
151 Performance was calculated as the percent of time after training in which the network output was
152 within criterion 0.5 from the target function. With the nonbinary, yet discrete signal, the network
153 achieved $85\% \pm 11\%$ performance compared to $93\% \pm 3\%$ performance with the binary $M(t)$
154 modulatory signal (95% confidence level). While these performance levels are not significantly
155 different, the larger standard error with the nonbinary signal indicates that with $M(t) = 0.1$ or 0.9 ,
156 working memory emergence and performance are less consistent per trial and potentially more
157 vulnerable to random network configurations, chaoticity, and noise (Figure 3). This is likely
158 because the baseline discrete signal permits weight change even if performance is significantly
159 poorer than average, leading to inappropriate weight change and failure to converge to the target
160 output.

161



162

163 Figure 3. Comparison of network output (blue) and target function (green)
164 modulatory signal (left panel) or a nonbinary, but discrete signal (right panel). The red arrows
165 represent the end of training.

166

167 3.2.2 Implementing an analog modulatory signal

168 In the brain, dopaminergic neurons fire with tonic activity, and their activity has been observed to
169 encode a sort of reward signal. More specifically, dopamine neurons have been recorded firing in
170 ways that represent a sort of “reward prediction error” [2]. This means roughly that upon
171 presentation of an unexpected reward, dopaminergic neurons will increase their firing rate above
172 baseline, while when a predicted reward is not observed, they will decrease their firing rate below
173 baseline. The model we investigated did not originally include an analog reward signal with
174 similar dynamics to dopaminergic neurons. Thus, we modified the model by changing the reward
175 signal to a function that approximated the behavior of the dopaminergic neurons.

176 In line with the original model, we thought of the average performance computed with a moving-
177 average filter as being analogous to the predicted reward. Thus, at any time, the current “reward
178 prediction error” would be given by the difference between the current performance P and the
179 average performance P_{avg} . We then used this difference as the independent variable in a linear
180 equation giving the reward signal. We set the intercept of this equation to be 0.1 to represent a
181 constant level of dopamine from tonically active dopaminergic neurons. We set the slope of the
182 linear function to be $(0.9-0.1) = 0.8$ to approximate the conditions imposed by the previous non-
183 binary discrete reward signal. Finally, a nonlinearity was added in the form of a hard cutoff at 0: If
184 the linear function produced values below 0, it was artificial set to 0. The combination of a strict
185 lower bound, a nonzero baseline signal, and dynamics depending upon a difference between a
186 stored and observed variable make this function a reasonable approximation to the observe
187 behavior of dopaminergic neurons.

188 The network was successfully able to learn the working memory task using the “dopamine-like”
189 reward signal (Figure 4). Interestingly, the network showed a slight increase in performance
190 during the testing period compared to the original model. The model trained with the “dopamine-
191 like” reward signal achieved a performance of $(96.2 \pm 0.99)\%$. This contrasts with the original
192 model which achieved a $(93 \pm 3)\%$ performance. Thus, for the modified model, the performance
193 increased and the variation in performance decreased.

194

195

196

197

198

199

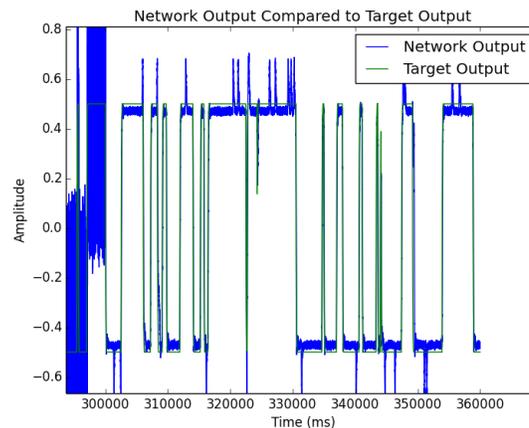
200

201

202

203

204

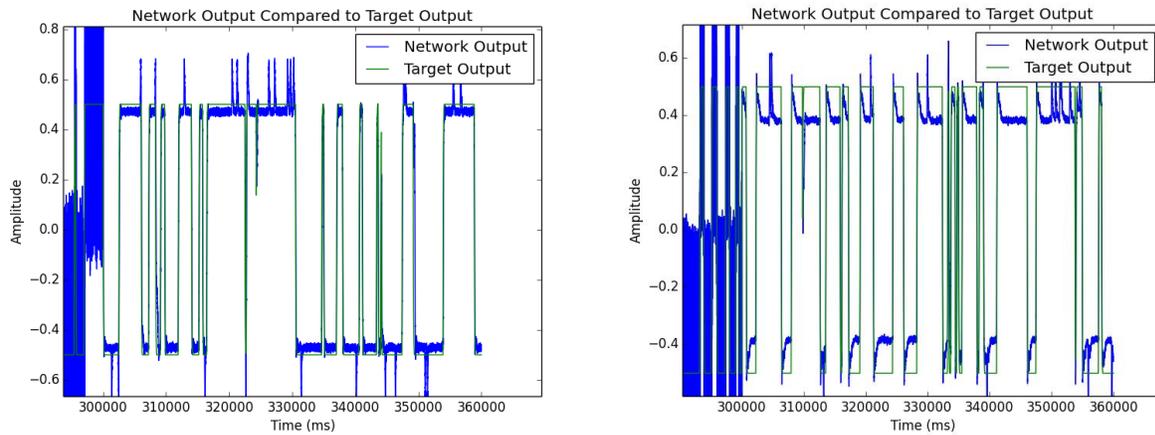


205 Figure 4. Comparison between the network output after training and the working memory target
206 function when using an analog “dopamine-like” modulatory signal

207

208 The analog reward signal we have described is not the only analog reward signal that could be
209 constructed. One obvious choice is the raw difference between the current performance and the
210 average (predicted) performance. For comparison, we used this raw difference to train a different
211 model to perform the working memory task. This new network was also able to successfully learn
212 the task, but a comparison of the network output of this raw performance difference network to the
213 output of the “dopamine-like” network show drastic differences. Namely, the raw performance

214 difference network consistently failed to follow the target function as precisely as the other
215 models. While the transitions were correct, the amplitudes and variances of the output signals were
216 very different and led to a poorer approximation. The figure below illustrates the differences seen
217 between the “dopamine-like” network and the raw performance difference network.



218
219 Figure 5. Comparison of working memory emergence with a “dopamine-like” modulatory signal
220 (left panel) and a raw performance difference signal (right panel)

221
222 One possible explanation for the apparent improvement in performance for the “dopamine-like”
223 reward signal could be based on the fact that a baseline reward signal implies that the weights are
224 changing if performance is not significantly worse than average. In this way, even when the
225 network performance dips slightly below average, the weights still change to explore new regions
226 of weight space. This contrasts with using the raw performance difference signal, in which the
227 weight change can be anti-Hebbian. This might imply that if training finds a set of weights that
228 lead to a relatively constant performance, then if the network tries to explore new regions of
229 weight space with poorer performance, anti-Hebbian learning will prevent the weights from
230 leaving that region. Thus, it is conceivable that using the raw performance difference signal
231 increases the probability of the network to fall into local minima. However, the constant weight
232 change of the “dopamine-like” reward signal allows the network to find more global minima so
233 long as performance is not significantly poorer than average, in which case the hard cutoff signal
234 of zero prevents further exploration.

235 236 3.3 Development of a delayed non-matching to sample task

237 In order to examine biological relevance of our model with modified modulatory signal, we
238 conceptualized a behavioral working memory test called delayed non-matching to sample task.
239 The test consists of three phases: sample, delay, and test phase. During the sample phase, the
240 rodent is presented with a sample stimulus (e.g. a left lever). After the rodent presses the sample
241 lever, a food pallet is presented at the opposite side of the chamber (delay phase). After a short
242 delay, the sample stimulus (i.e. left lever) is shown again along with a novel alternative (i.e. right
243 lever). In this paradigm, the rodent is rewarded with food pallet for selecting the novel stimulus
244 (i.e. right lever). This delayed non-matching to sample task is a good behavioral assessment for
245 working memory emergence as it requires the animal to hold information about the sample
246 stimulus in the online workspace throughout the delay phase in order to correctly select the novel
247 stimulus during test phase (Figure 6). In our model, the sample stimulus and the novel stimulus
248 can be thought of as input 1 and input 2 respectively. During the sample phase, the network is
249 presented only with input 1, which has to be retained in the network’s working memory
250 throughout the delay phase. In the test phase, the network is presented with both input 1 and input

251 2. In order for network output to match the target function, the network has to correctly reject
252 input 1 and select input 2 during the test phase.

253

254

255

256

257

258

259

260

261

262

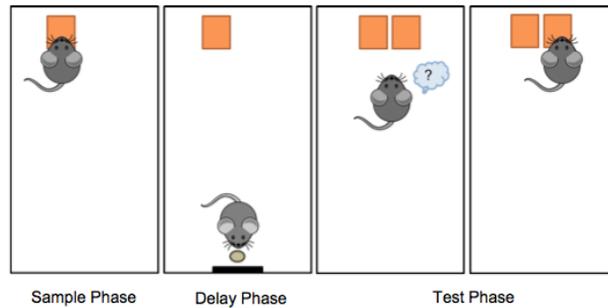


Figure 6. Schematic of delayed non-matching to sample task

263

264 In the delayed non-matching to sample task, we first tested working memory using both the binary
265 and nonbinary, yet discrete modulatory signals $M(t)$. With these signals, performance levels were
266 $71\% \pm 5\%$ and $72\% \pm 10\%$, respectively (confidence level 95%). These performance levels were
267 significantly worse than those with either discrete modulatory signal in the original working
268 memory task ($p < 0.001$). Thus, the network performed more poorly on the delayed non-matching
269 to sample task than on the original task. In particular, analysis of individual runs indicates that the
270 network was able to achieve correct output when one input stream was nonzero, but not when both
271 input streams were nonzero. The network was not able to retain information about the original
272 stimulus and switch its output when presented with the original stimulus and the novel stimulus. It
273 is possible that different types of network architecture or non-Hebbian learning rules could lead to
improved performance with the delayed non-matching to sample task.

274

275

276

277

278

279

280

281

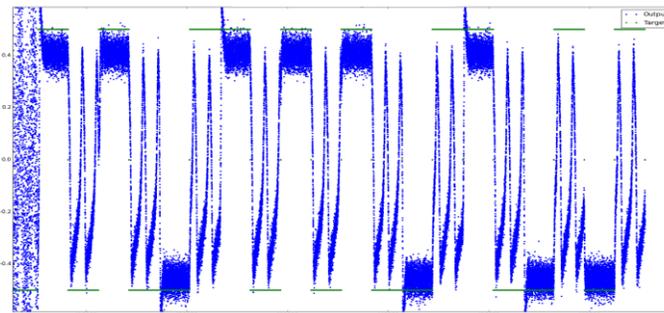


Figure 7. Failure to converge to target function in the delay non-matching to sample task

283

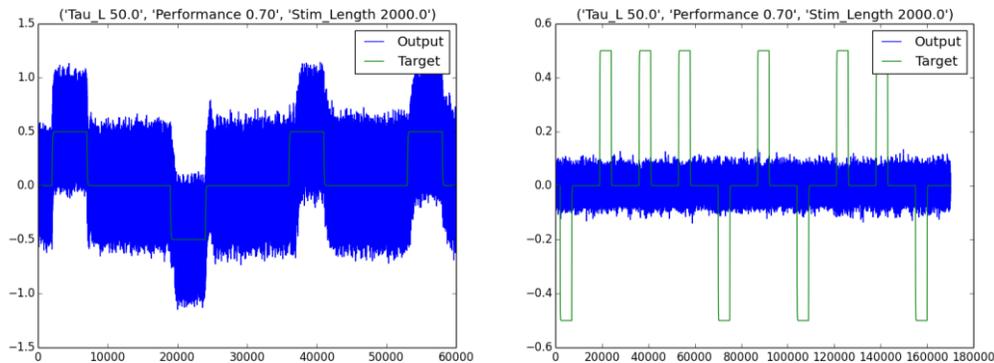
284

3.3.1 Implementing a delayed non-matching to sample task with a “go signal”

285

286 In an attempt to get the network to learn some version of the delayed non-matching to sample task,
287 we added a third input stream to serve as a “go signal.” The rationale for its inclusion was the
288 possibility that since in the original implementation of this new task, the individual input streams
289 require different outputs at different times, and that the network does not implement an explicit
290 temporal difference learning rule, the network was unable to differentiate the meanings of the
291 original input signals. We redesigned the task to present one pulse on one of the input streams,
292 followed by a pulse on the “go signal” stream. Upon this pulse, we directed the network to choose
293 the opposite value compared to the input stream. Thus, for every trial of this task, either original
input stream requires only one output at any time. A binary modulatory signal was used.

294 Unfortunately, the network performed even worse on this task. The figure below shows the
295 network output during and after training compared to the target function.
296



297

298 Figure 8. Comparison of network output and target function before (left panel) and after (right
299 panel) training with a go-signal delay non-matching to sample task
300

301 Following 500 seconds of training, the network showed what appeared to be random noise
302 throughout the testing period. No evidence was seen for any state transitions that would be
303 indicative of the network trying to accomplish the task. Future work will have to investigate
304 whether there is a fundamental reason the network was unable to learn this task.

305

306 4 Conclusions

307 This study examines working memory emergence using LSM framework. While several other
308 studies trained their recurrent neural networks using supervised paradigms, the current study used
309 unsupervised, reward-modulated Hebbian learning to examine working memory emergence in
310 recurrent neural networks. We first reproduced a model proposed in [1] and showed matching
311 target output in recurrent neural networks via reward-modulated Hebbian learning.

312 In addition, we replaced the binary modulatory signal used in to the original model [1] with a
313 nonbinary, but discrete signal that represented the dynamics of the networks more accurately. We
314 further modified the reward signal to reflect dopaminergic neurons' "reward prediction error,"
315 where dopaminergic neurons increase their firing rate above baseline when an unexpected reward
316 is presented, and decrease their firing rate below baseline when a predicted reward is not observed.
317 To incorporate this behavior of dopaminergic neurons into our model, we implemented an analog
318 modulatory signal where, at any time, the current "reward prediction error" was given by the
319 difference between the current performance and the average performance on the task, though with
320 a lower limit of zero. We demonstrated that the our modified model with the "dopamine-like"
321 reward signal, compared to the discrete signals, led to increased performance and decreased
322 variation in performance on the original working memory task.

323 Finally, in order to further examine working memory emergence, we designed a behavioral
324 working memory test (delayed non-matching to sample task) and trained the network to learn it.
325 This task required the network to hold information regarding a sample stimulus (input 1) through
326 the delay period and select a novel stimulus (input 2) when both novel and sample stimuli are
327 presented during the test phase. However, we found that the neural network performed more
328 poorly on the delayed non-matching to sample task than on the original task: the network could
329 not achieve output matching the target function when two input streams were presented (i.e.,
330 nonzero). Future studies will be aimed at determining what network architectures, modulatory
331 signals, or learning rules will facilitate the learning of the non-matching to sample task.

332 **References**

- 333 [1] Hoerzer, G.M., Legenstein, R., & Maass, W. (2014) Emergence of complex computational
334 structures from chaotic neural networks through reward-modulated Hebbian learning. *Cerebral*
335 *Cortex* **3**: 677-690.
- 336 [2] Lak, A., Stauffer, W. R., & Schultz, W. (2014) Dopamine prediction error responses integrate
337 subjective value from different reward dimensions. *PNAS* **111**(6): 2343-2348.
- 338 [3] Sussillo, D., and Abbott, L.F. (2009). Generating coherent patterns of activity from chaotic
339 neural networks. *Neuron* **63**: 544-557.