

1 Flexible gating between subspaces by a disinhibitory motif: a
2 neural network model of internally guided task switching

3 Yue Liu*, Xiao-Jing Wang*
4 New York University

5 Abstract

Behavioral flexibility relies on the brain's ability to switch rapidly between multiple tasks, even when the task rule is not explicitly cued but must be inferred through trial and error. The underlying neural circuit mechanism remains poorly understood. We investigated recurrent neural networks (RNNs) trained to perform an analog of the classic Wisconsin Card Sorting Test. The networks consist of two modules responsible for rule representation and sensorimotor mapping, respectively, where each module is comprised of a circuit with excitatory neurons and three major types of inhibitory neurons. We found that rule representation by self-sustained persistent activity across trials, error monitoring and gated sensorimotor mapping emerged from training. Systematic dissection of trained RNNs revealed a detailed circuit mechanism that is consistent across networks trained with different hyperparameters. The networks' dynamical trajectories for different rules reside in separate subspaces of population activity; they become virtually identical and performance was reduced to chance level when dendrite-targeting somatostatin-expressing interneurons were silenced, demonstrating that rule-based gating critically depends on the disinhibitory motif.

7 **Introduction**

8 A signature of cognitive flexibility is the ability to adapt to a changing task demand.
9 Oftentimes, the relevant task is not explicitly instructed, but needs to be inferred from
10 previous experiences. In laboratory studies, this behavioral flexibility is termed un-cued
11 task switching. A classic task to evaluate this ability is the Wisconsin Card Sorting Test

*Corresponding author. Email: y14317@nyu.edu, xjwang@nyu.edu

12 (WCST) [1]. During this task, subjects are presented with an array of cards, each with
13 multiple features, and should respond based on the relevant feature dimension (i.e. the task
14 rule) that changes across trials. Crucially, subjects are not instructed on when the rule
15 changes, but must infer the currently relevant rule based on the outcome of previous trials.
16 Intact performance on un-cued task switching depends on higher-order cortical areas such
17 as the prefrontal cortex (PFC) [2, 3, 4, 5, 6], which has been proposed to represent the task
18 rule and modulate the activity of other cortical areas along the sensorimotor pathway [7].

19 Four essential neural computations must be implemented by the neural circuitry un-
20 derlying un-cued task switching. First, it should maintain an internal representation of
21 the task rule across multiple trials when the rule is unchanged. Second, soon after the
22 rule switches, the animal will inevitably make errors and receive negative feedback, since
23 the switches are un-cued. This negative feedback should induce an update to the internal
24 representation of the task rule. Third, the neural signal about the task rule should be
25 communicated to the brain regions responsible for sensory processing and action selection.
26 Fourth, this rule signal should be integrated with the incoming sensory stimulus to produce
27 the correct action.

28 Prior work has identified neural correlates of cognitive variables presumed to underlie
29 these computations including rule [8], feedback [8, 9, 10] and conjunctive codes for sensory,
30 rule, and motor information [11]. In addition, different types of inhibitory neurons are
31 known to play different functional roles in neural computation: while parvalbumin (PV)-
32 expressing interneurons are suggested to underlie feedforward inhibition [12], interneurons
33 that express somatostatin (SST) and vasoactive intestinal peptide (VIP) have been pro-
34 posed to mediate top-down control [13, 14, 15, 16]. In particular, SST and VIP neurons
35 form a disinhibitory motif [17, 18, 19] that has been hypothesized to instantiate a gating
36 mechanism for flexible routing of information in the brain [20]. However, there is currently a
37 lack of mechanistic understanding of how these neural representations and cell-type-specific
38 mechanisms work together to accomplish un-cued task switching.

39 To this end, we used computational modeling to formalize and discover mechanistic
40 hypotheses. In particular, we used tools from machine learning to train a collection of
41 biologically informed recurrent neural networks (RNNs) to perform an analog of the WCST
42 used in monkeys [21, 8, 9]. Training RNN [22] does not presume a particular circuit solution,
43 enabling us to explore potential mechanisms. For this purpose, it is crucial that the model
44 is biologically plausible. To that end, each RNN was set up to have two modules: a “PFC”
45 module for rule maintenance and switching and a “sensorimotor” module for executing the
46 sensorimotor transformation conditioned on the rule. To explore the potential functions
47 of different neuronal types in this task, each module of our network consists of excitatory
48 neurons with somatic and dendritic compartments as well as PV, SST and VIP inhibitory

49 neurons, where the connectivity between cell types is constrained by experimental data
50 (Methods).

51 After training, we performed extensive dissection of the trained models to reveal
52 that close interplay between local and across-area processing was essential for solving the
53 WCST. First, we found that abstract cognitive variables were distinctly represented in the
54 PFC module. In particular, two subpopulations of excitatory neurons emerge in the PFC
55 module - one encodes the task rule and the other shows mixed-selectivity that nonlinearly
56 depends on rule and negative feedback. Notably, neurons with similar response profiles
57 have been reported in neurophysiological recordings of monkeys performing the same task
58 [8, 9]. Second, we identified interesting structures in the local connectivity between differ-
59 ent neuronal assemblies within the PFC module, which enabled us to compress the high-
60 dimensional PFC module down to a low-dimensional simplified network. Importantly, the
61 neural mechanism for maintaining and switching rule representation is readily interpretable
62 in the simplified network. Third, we found that the rule information in the PFC module is
63 communicated to the sensorimotor module via structured long-range connectivity patterns
64 along the monosynaptic excitatory pathway, the di-synaptic pathway that involves PV neu-
65 rons, as well as the trisynaptic disinhibitory pathway that involves SST and VIP neurons.
66 In addition, different dendritic branches of the same excitatory neuron in the sensorimotor
67 module can be differentially modulated by the task rule depending on the sparsity of the lo-
68 cal connections from the dendrite-targeting SST inhibitory neurons. Fourth, single neurons
69 in the sensorimotor module show nonlinear mixed selectivity to stimulus, rule and response,
70 which crucially depends on the activity of the SST neurons. On the population level, the
71 neural trajectories for the sensorimotor neurons during different task rules occupy nearly
72 orthogonal subspaces, which is disrupted by silencing the SST neurons. Lastly, we found
73 structured patterns of input and output connections for the sensorimotor module, which
74 enables appropriate rule-dependent action selection. These results are consistent across
75 dozens of trained RNNs with different types of dendritic nonlinearities and initializations,
76 therefore pointing to a common neural circuit mechanism underlying the WCST.

77 Results

78 **Training modular recurrent neural networks with different types of inhibitory** 79 **neurons**

80 We trained a collection of modular RNNs to perform the WCST. Each RNN consists
81 of two modules: the “PFC” module receives an input about the outcome of the previous
82 trial, and was trained to output the current rule; the “sensorimotor” module receives the
83 sensory input and was trained to generate the correct choice (Figure 1b). The inputs
84 and outputs were represented by binary vectors (Figure 1b, Methods) Each module was

85 endowed with excitatory neurons with somatic and two dendritic compartments, as well as
86 three major types of genetically-defined inhibitory neurons (PV, SST and VIP). Different
87 types of neurons have different inward and outward connectivity patterns constrained by
88 experimental data in a binary fashion (Methods, Figure 1b). The somata of all neurons were
89 modeled as standard leaky units with a rectified linear activation function. The activation
90 of the dendritic compartments, which can be viewed as a proxy for the dendritic voltage, is a
91 nonlinear sigmoidal function of the excitatory and inhibitory inputs they receive (Methods).
92 The specific form of the nonlinearity is inspired by experiments showing that inhibition acts
93 subtractively or divisively on the dendritic nonlinearity function depending on its relative
94 location to the excitation along the dendritic branch [23]. Therefore, we trained a collection
95 of RNNs, each with either subtractive or divisive dendritic nonlinearity, to explore the effect
96 of dendritic nonlinearity on the network function.

97 The task we trained the network on is a WCST variant used in monkey experiments
98 [21, 8, 9, 6] (Figure 1a). During each trial, a reference card with a particular color and
99 shape is presented on the screen for 500 ms, after which three test cards appear around
100 the reference card for another 500 ms. Each card can have one of the two colors (red or
101 blue) and one of the two shapes (square or triangle). A choice should be made for the
102 location that contains the test card that has the same relevant feature (color or shape) as
103 the reference card, after which the outcome of the trial is given, followed by an inter-trial
104 interval. The relevant feature to focus on, or the task rule, changes randomly every few
105 trials. Critically, the rule changes were not cued, requiring the network to memorize the rule
106 of the last trial using its own dynamics. Therefore, the network dynamics should be carried
107 over between consecutive trials, rather than reset at the end of each trial as has been done
108 traditionally [24, 25]. To this end, the network operated continuously across multiple trials
109 during training, and the loss function was aggregated across multiple trials (Methods). We
110 use supervised learning to adjust the strength of all the connections (input, recurrent and
111 output) by minimizing the mean squared error between the output of both modules and
112 the desired output (rule for the PFC module and response for the sensorimotor module).
113 Notably, only the connections between certain cell types are non-zero and can be modified.
114 This is achieved using a mask matrix, similar to [26] (Methods).

115 After training converged, we tested the models by running them continuously across
116 100 trials of WCST with 10 rule switches at randomly chosen trials. The networks made a
117 single error after each rule switch, and were able quickly switch to the new rule and main-
118 tain good performance (Figure 1c, d). Correspondingly, single neurons from both modules
119 exhibited rule-modulated persistent activity that lasted several trials (Supplementary Fig-
120 ure 1).

121 Our networks can reliably maintain good performance after a single correct trial in the

122 new rule, which matches the behavior of monkeys in some previous studies (e.g. Figure 1e).
123 However, the performance of monkeys during this task showed substantial variability across
124 different studies as well as different sessions with a study. The number of error trials that
125 monkeys take to switch to the new rule ranges from one trial to tens of trials [21, 8, 9, 6]. One
126 reason is that performance typically reaches a certain criterion (e.g. 80% correct) but not
127 perfect accuracy before rule switching, therefore an error signal could mean an erroneous
128 sensori-motor decision rather than rule change. Indeed, when training of our model is
129 stopped at 80% rather than 100% accuracy, the resulting network shows gradual switching
130 (Figure 1f, Methods). This point will be addressed further in the Discussion section. In the
131 following sections, we will “open the black box” to understand the mechanism the networks
132 used to perform the WCST.

133 **Two rule attractor states in the PFC module maintained by interactions between** 134 **modules**

135 We first dissected the PFC module, which was trained to represent the rule. Since
136 there are two rules in the WCST task we used, we expected that the PFC module might have
137 two attractor states corresponding to the two rules. Therefore, we examined the attractor
138 structure in the dynamical landscape of the PFC module by initializing the network at
139 states chosen randomly from the trial, and evolving the network autonomously (without any
140 input) for 500 time steps (which equals 5 seconds in real time). Then, the dynamics of the
141 PFC module during this evolution was visualized by applying principal component analysis
142 to the population activity. The PFC population activity settled into one of two different
143 attractor states depending on the rule that the initial state belongs to (Supplementary
144 Figure 2a). Therefore, there are two attractors in the dynamical landscape of the PFC
145 module, corresponding to the two rules.

146 Historically, persistent neural activity corresponding to attractor states were first
147 discovered in the PFC [28, 29, 30, 31]. However, more recent experiments found persistent
148 neural activity in multiple brain regions, suggesting that long-range connections between
149 brain regions may be essential for generating persistent activity [32, 33, 34, 35]. Inspired
150 by these findings, we wondered if the PFC module in our network could support the two
151 rule attractor states by itself, or that the long-range connections between the PFC and the
152 sensorimotor module are necessary to support them. To this end, we lesioned the inter-
153 modular connections in the trained networks and repeated the simulation. Interestingly,
154 we found that for the majority of the trained networks (42 out of 52 for the fast switching
155 networks and 3 out of 3 for the slow switching networks), their PFC activity settled into a
156 trivial fixed point corresponding to an inactive state (Supplementary Figure 2b, c). This
157 result shows that the two rule attractor states in these networks are dependent on the

DISSECTING MODULAR RNNs TRAINED TO PERFORM A WCST ANALOG 6

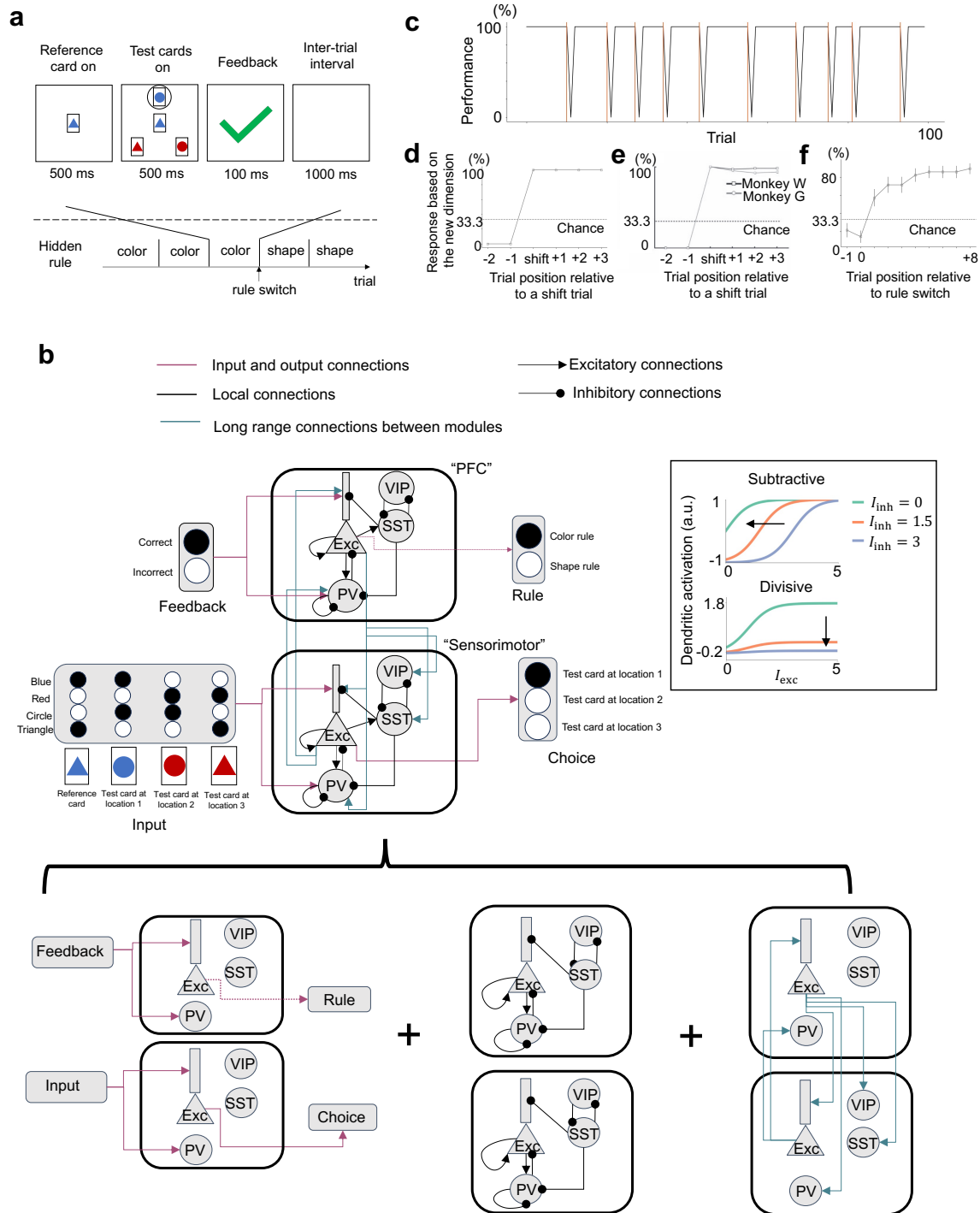


Figure 1: (Caption next page.)

Figure 1: (Previous page.) **Model setup and performance.**

a. The schematic of the WCST task. Subjects are required to choose the card that matches the reference card at the center in either shape or color, depending on a hidden rule that switches after a number of trials.

b. The RNN contains a “PFC” module and a “sensorimotor” module. The PFC module receives an input about the feedback of the previous trial, and was trained to produce the current rule. The sensorimotor module receives the sensory input and was trained to produce the correct choice. The input to the PFC module about the feedback is represented by a two-dimensional binary vector. The input to the sensorimotor module represents the features of the cards on the screen. Each card is represented by a four-dimensional binary vector, where the two non-zero entries represent the color and shape of the card. The target output of the PFC module about the correct rule is represented by a two-dimensional binary vector. The target output of the sensorimotor module about the correct choice is represented by a three-dimensional binary vector. Each module is endowed with excitatory neurons and three types of inhibitory neurons: PV, SST and VIP. The cell-type-specific connectivity is constrained by experimental data (Methods). Bottom panel shows the decomposition of the model architecture into the input and output connectivity (left, magenta). The dashed line from PFC to rule represents the fact that the PFC module was trained to represent the rule but there are no explicit rule output neurons in the model), the local recurrent connectivity (middle, black) and inter-modular connectivity (green, right). All connections were trained. Each excitatory neuron is modeled with a somatic and two dendritic compartments. Inset shows for the two types of dendritic nonlinearities used the relationship between the excitatory input onto the dendrite and the dendritic activity for different levels of inhibitory inputs.

c. The performance of the model during testing, for an example network. The network made one error after each rule switch (red vertical lines) and quickly recovered its performance.

d. Performance as a function of trial position relative to the first correct trial after rule change, or the “shift” trial, for the same example network as in **c**.

e. The performance of two monkeys as a function of trial position relative to the shift trial. Figure adapted from Ref. [27]

f. The performance of an example model where training was stopped before it reached perfect performance. This model exhibit more gradual switching between rules.

158 interactions between the PFC and the sensorimotor modules.

159 **Two emergent subpopulations of excitatory neurons in the PFC module**

160 For the PFC module to keep track of the current rule in effect, the module should
161 stay in the same rule attractor state after receiving positive feedback, but transition to the
162 other rule attractor state after receiving negative feedback. We reasoned that this network
163 function might be mediated by single neurons that are modulated by the task rule and
164 negative feedback, respectively. Therefore, we set out to look for these single neurons.

165 In the PFC module of the trained networks, there are indeed neurons whose activity
166 is modulated by the task rule in a sustained fashion (example neurons in Supplementary
167 Figure 1 and Figure 2a, top). In contrast, there are also neurons that show transient
168 activity only after negative feedback. Furthermore, this activity is also rule-dependent. In
169 other words, their activity is conjunctively modulated by negative feedback and the task
170 rule (example neurons in Supplementary Figure 1, red traces and Figure 2a, bottom). We
171 termed these two classes of neurons “rule neurons” and “conjunctive error x rule neurons”
172 respectively.

173 We identified all the rule neurons and conjunctive error x rule neurons in the PFC
174 module using a single neuron selectivity measure (see Methods for details). The two classes
175 of neurons are clearly separable on the two-dimensional plane in Figure 2c, where the x axis
176 is the input weight for negative feedback, and the y axis is the rule modulation, which is the
177 difference in the mean activity between the two rules (for trials following a correct trial).
178 As shown in Figure 2c, rule neurons receive negligible input about negative feedback, and
179 many of them have activity modulated by rule. On the other hand, conjunctive error x rule
180 neurons receive a substantial amount of input about negative feedback, yet their activity is
181 minimally modulated by rule on trials following a correct trial (Figure 2b). This pattern
182 was preserved when aggregating across trained networks (Figure 2c and Supplementary
183 Figure 3). Interestingly, neurons with similar tuning profiles have been reported in the
184 PFC and posterior parietal cortex of macaque monkeys performing the same WCST analog
185 [8, 9].

186 Across different cell types in the PFC module, on average 23.1% of excitatory neurons,
187 57.3% of PV neurons and 38.1% of SST neurons were classified as rule neurons in each model.
188 Compared to excitatory neurons, a much smaller fraction of inhibitory neurons in the PFC
189 were classified as conjunctive error x rule neurons. On average, 22.9% excitatory neurons
190 were conjunctive error x rule neurons in each model, compared with 11.5% PV neurons
191 and 5.2% SST neurons. Therefore, we focus only on the excitatory conjunctive error x rule
192 neurons in the analysis below.

193 We also performed the same analysis on the trained networks that switch rules more

194 slowly (e.g. Figure 1f). In those networks there is not a clear separation between the two
195 subpopulations of excitatory neurons in the PFC module (data not shown).

196 **Maintaining and switching rule states via structured connectivity patterns be-**
197 **tween subpopulations of neurons within the PFC module**

198 Given the existence of rule neurons and conjunctive error x rule neurons, what is the
199 connectivity between them that enables the PFC module to stay in the same rule attractor
200 state when receiving correct feedback, and switch to the other rule attractor state when
201 receiving negative feedback?

202 To this end, we examined the connectivity between different subpopulations of neu-
203 rons in the PFC module explicitly, by computing the mean connection strength between
204 each pair of subpopulations. This analysis reveals that the excitatory rule neurons and
205 PV rule neurons form a classic winner-take-all network architecture [36] with selective in-
206 hibitory populations [37, 38], where excitatory neurons preferring the same rule are more
207 strongly connected, and they also more strongly project to PV neurons preferring the same
208 rule (Figure 3a). On the other hand, PV neurons project more strongly to both excitatory
209 neurons and other PV neurons with the opposite rule preference (Figure 3a). This winner-
210 take-all network motif together with the excitatory drive from the sensorimotor module
211 (Supplementary Figure 2) is able to sustain one of the two attractor states.

212 Next, how are the rule neurons connected with the conjunctive error x rule neurons
213 such that the sub-network formed by rule neurons can switch from one attractor to the other
214 in the presence of the negative feedback input? Using the same method, we found that the
215 connectivity between the rule neurons and the conjunctive error x rule neurons exhibited
216 an interesting structure: the excitatory rule neurons more strongly target the conjunctive
217 error x rule neurons that prefer the opposite rule; the PV rule neurons more strongly target
218 conjunctive error x rule neurons that prefer the same rule (Figure 3b, top two panels). On
219 the other hand, the conjunctive error x rule neurons more strongly target the excitatory
220 and PV rule neurons that prefer the same rule (Figure 3b, bottom two panels).

221 This connectivity structure gives rise to a simple circuit diagram of the PFC module
222 (Figure 3c), which leads to an intuitive explanation of the circuit mechanism underlying the
223 switching of rule attractor state. For example, suppose the network is in the attractor state
224 corresponding to color rule, and has just received a negative feedback and is about to switch
225 to the attractor corresponding to the shape rule (Figure 3e, left). As shown in Figure 2b-c,
226 the input current that represents the negative feedback mainly targets the conjunctive error
227 x rule neurons. In addition, since the network is in the color rule state, the excitatory and
228 PV neurons that prefer the color rule are more active than those that prefer the shape
229 rule. According to Figure 3b (top two panels), the excitatory neurons that prefer the color

DISSECTING MODULAR RNNs TRAINED TO PERFORM A WCST ANALOG 10

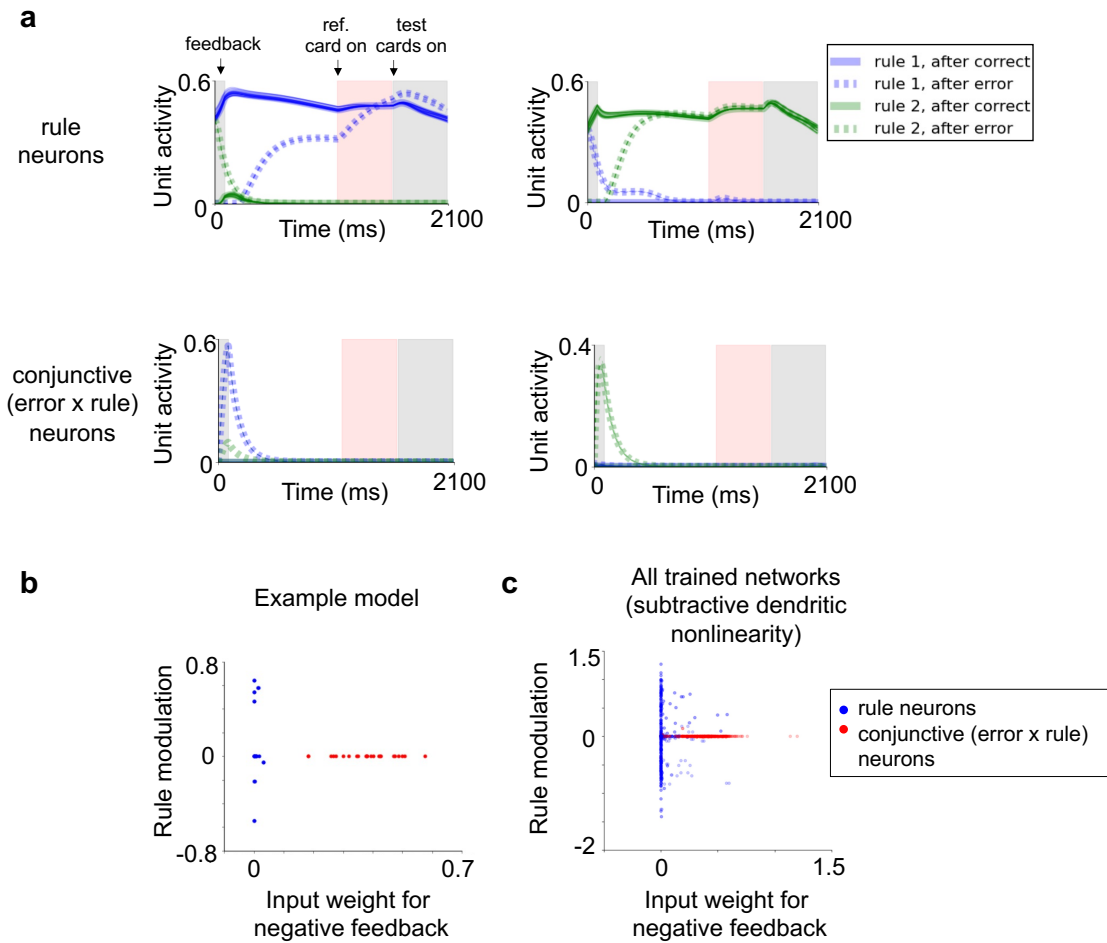


Figure 2: Emergence of two subpopulations of excitatory neurons in the PFC module after training.

a. Two example rule neurons (top) and conjunctive error x rule neurons (bottom). The solid traces represent the mean activity across trials that follows a correct trial, when those trials belong to color rule (blue) or shape rule (green) blocks. The dashed traces represent the mean activity after error trials, when those trials belong to color rule (blue) or shape rule (green) blocks. We use rule 1 and color rule, as well as rule 2 and shape rule interchangeably hereafter.

b. Rule neurons and conjunctive neurons are separable. The x axis represents the input weight for negative feedback, and the y axis is the difference between the mean activity over color rule trials and shape rule trials (for trials following a correct trial). As shown, the rule neurons (blue points) receive little input about negative feedback, but their activity is modulated by rule; The conjunctive error x rule neurons (red points) receive substantial input about negative feedback, but their activity is not modulated by rule (during trials following a correct trial).

c. The trend in **b** is preserved across a collection of trained networks. Here the results are shown for networks with subtractive dendritic nonlinearity. Networks with divisive dendritic nonlinearity show a similar pattern (Supplementary Figure 3).

230 rule strongly excite the error x shape rule neurons, and the PV neurons that prefer the
231 color rule strongly inhibit the error x color rule neurons. Therefore, the error x shape rule
232 neurons receive stronger total input than the error x color rule neurons, and will be more
233 active (Figure 3e, middle). Their activation will in turn excite the excitatory neurons and
234 PV neurons that prefer the shape rule (Figure 3b, bottom two panels). Finally, due to the
235 winner-take-all connectivity between the rule populations (Figure 3a), the excitatory and
236 PV neurons that prefer the color rule will be suppressed, and the network will transition to
237 the attractor state for the shape rule (Figure 3e, right).

238 It is worth noting that the same mechanism can also trigger a transition in the opposite
239 direction (from shape rule to color rule) in the presence of the same negative feedback
240 signal. This is enabled by the biased connections between the rule and conjunctive error x
241 rule populations.

242 Is the simplified circuit diagram (Figure 3c) consistent across trained networks, or
243 different trained networks use different solutions? To examine this question, we computed
244 a “connectivity bias” measure between each pair of populations for each trained network.
245 This measure is greater than zero if the connectivity structure between a pair of populations
246 is closer to the one in the simplified circuit diagram in Figure 3c than to the opposite
247 (see Methods for details). Across trained networks, we found that the connectivity biases
248 were mostly greater than zero (Figure 3d), indicating that the same circuit motif for rule
249 maintenance and switching underlies the PFC module across different trained networks.

250 A similar circuit architecture exists between the excitatory neurons and the SST neu-
251 rons in the PFC module (Supplementary Figure 5), where SST neurons receive stronger
252 excitatory input from the conjunctive error x rule neurons that prefer the same rule, and
253 also more strongly inhibit the error x rule neurons that prefer the same rule. In addi-
254 tion, they form a winner-take-all connectivity with the rule excitatory neurons by receiving
255 stronger projections from the rule neurons that prefer the same rule and projecting back
256 more strongly to the rule neurons that prefer the opposite rule. Therefore, they contribute
257 to rule maintenance and switching in a similar way as the PV neurons.

258 When the same analysis was performed on the slow-switching networks (e.g. Fig-
259 ure 1f), we found that although the rule neurons in these networks also form a winner-take-
260 all connectivity structure, the connectivity between error x rule neurons and rule neurons
261 does not exhibit the same structure as in the fast-switching models (data not shown). There-
262 fore, the slow-switching networks have a similar sub-network that encodes the rule, but a
263 poorly organized sub-network between the error x rule and rule neurons, which may explain
264 why switching the rule takes more trial-and-error in these networks.

DISSECTING MODULAR RNNs TRAINED TO PERFORM A WCST ANALOG 12

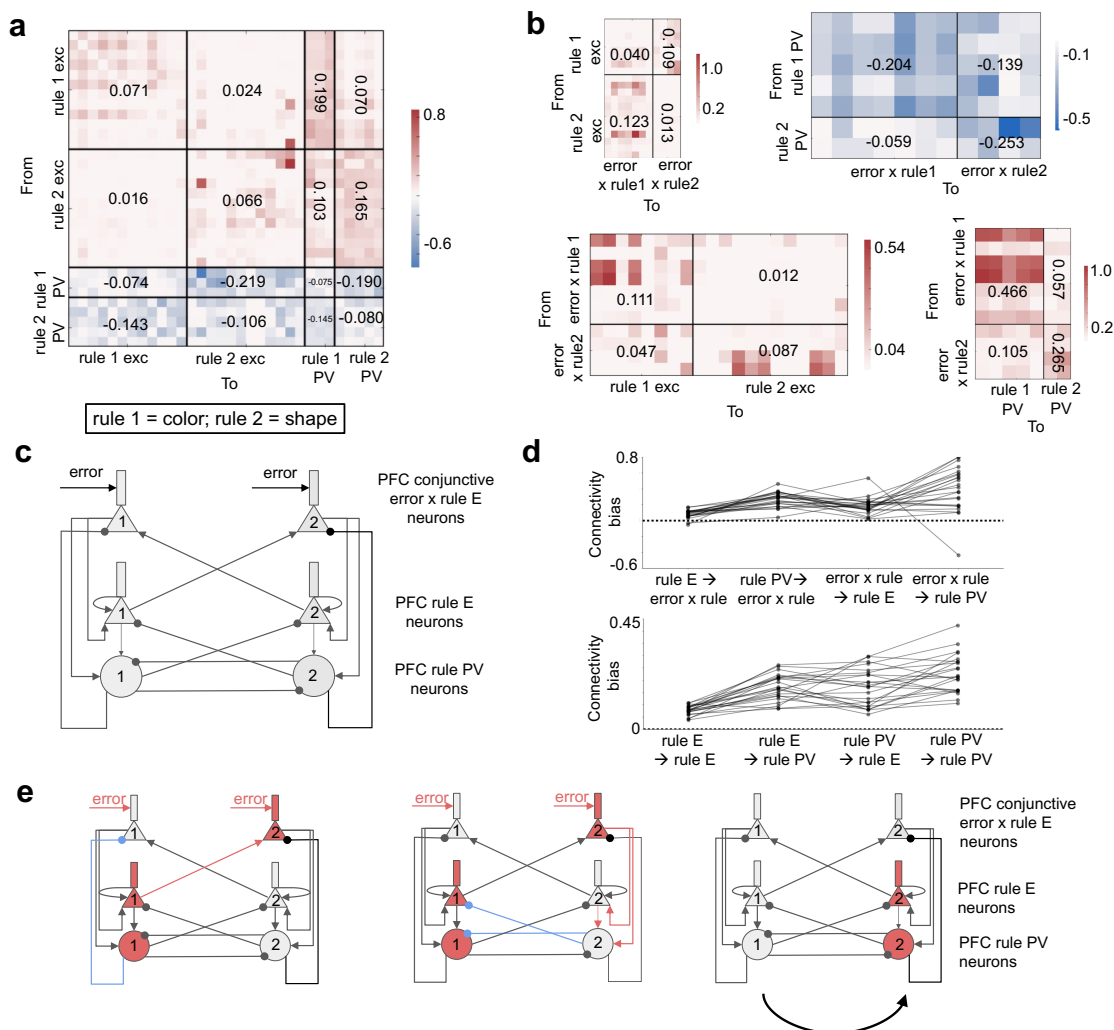


Figure 3: (Caption next page.)

265 **Top-down propagation of the rule information through structured long-range**
 266 **connections**

267 Given that the PFC module can successfully maintain and update the rule representation,
 268 how does it use the rule representation to reconfigure the sensorimotor mapping?
 269 First, we found that neurons in the sensorimotor module were tuned to rule (Supplementary
 270 Figure 6a), since they receive top-down input from the rule neurons in the PFC module. The
 271 PFC module exerts top-down control through three pathways: the monosynaptic pathway
 272 from the excitatory neurons in the PFC module to the excitatory neurons in the sensori-
 273 motor module, the tri-synaptic pathway that goes through the VIP and SST neurons in the
 274 sensorimotor module, and the di-synaptic pathway mediated by the PV neurons in the
 275 sensorimotor module (Figure 1b). We found that there are structured connectivity pat-

Figure 3: (Previous page.) **An emergent circuit wiring diagram in the PFC module enables un-cued switching between rule attractor states.**

a. The connectivity matrix between different populations of rule neurons, for an example model. Text indicates the mean connection strength between two populations. The excitatory rule neurons project more strongly to, and receive more input from, neurons with the same preferred rule. The PV rule neurons project more strongly to and receive more input from neurons with the opposite rule preference. As a result, rule neurons form a classic winner-take-all connectivity with selective inhibitory populations that maintain the two rule attractor state.

b. The connectivity between rule neurons and conjunctive error x rule neurons, for an example model. Top left: excitatory rule neurons project more strongly to the conjunctive error x rule neurons that prefer the opposite rule; Top right: PV rule neurons project more strongly to conjunctive error x rule neurons that prefer the same rule; Bottom left: conjunctive error x rule neurons project more strongly to the excitatory rule neurons that prefer the same rule; Bottom right: conjunctive error x rule neurons project more strongly to the PV rule neurons that prefer the same rule.

c. The simplified circuit diagram between rule neurons and conjunctive neurons based on the result of **b**. The weaker connections are ignored. Rule 1 represents the color rule and rule 2 represents the shape rule

d. A connectivity bias was computed to describe the extent to which the connectivity pattern between each pair of subpopulations conform to the simplified diagram in **c**. A value greater than 0 indicates the connectivity structure is more similar to that in **c** than to the opposite. The connectivity biases across all trained models are mostly above 0, both for the connection among rule neurons (top) and the connection between rule neurons and conjunctive error x rule neurons (bottom). Here the results are shown for networks with subtractive dendritic nonlinearity. Networks with divisive dendritic nonlinearity show similar result (Supplementary Figure 4).

e. A schematic showing how the simplified circuit can switch from the rule 1 attractor state to the rule 2 attractor state after receiving the input about negative feedback. The conjunctive error x rule 2 neurons receive excitation from the currently-active rule 1 excitatory neurons (red arrow, left panel), and the conjunctive error x rule 1 neurons receive inhibition from the currently-active rule 1 PV neurons (blue arrow, left panel). This makes conjunctive error x rule 2 neurons more active than the conjunctive error x rule 1 neurons, even though the negative feedback input targets both error x rule 1 and error x rule 2 populations (left panel). The conjunctive error x rule 2 neurons then excite the rule 2 excitatory and PV neurons (red arrows, middle), which suppress the rule 1 excitatory and PV neurons due to the winner-take-all connectivity (blue arrows, middle) and eventually become more active (right).

276 terns along all three pathways. Along the monosynaptic pathway, excitatory rule neurons
277 in the PFC module preferentially send long-range projections to the excitatory neurons
278 in the sensorimotor module that prefer the same rule (Figure 4a). Along the tri-synaptic
279 pathway, PFC excitatory rule neurons also send long-range projections to the SST and VIP
280 interneurons in the sensorimotor module that prefer the same rule (Figure 4b-c). The SST

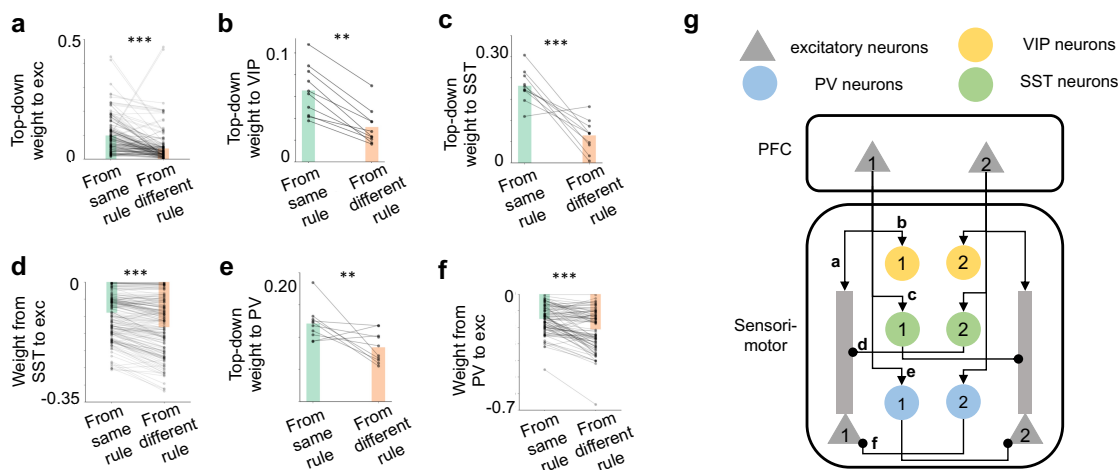


Figure 4: (Caption next page.)

281 neurons in turn send stronger inhibitory connections to the dendrite of the local excita-
 282 tory neurons that prefer the opposite rule (Figure 4d). Along the di-synaptic pathway, the
 283 PV neurons are also more strongly targeted by PFC excitatory rule neurons that prefer the
 284 same rule (Figure 4e), and they inhibit local excitatory neurons that prefer the opposite rule
 285 (Figure 4f). These trends are preserved across trained networks (Supplementary Figure 7).
 286 Therefore, rule information is communicated to the sensorimotor module synergistically via
 287 the mono-synaptic excitatory pathway, the tri-synaptic pathway that involves the SST and
 288 VIP neurons, as well as the di-synaptic pathway that involves the PV neurons, as illustrated
 289 in Figure 4g.

290 **Structured input and output connections of the sensorimotor module enable**
 291 **rule-dependent action selection**

292 How does the sensorimotor module implement the sensorimotor transformation (from
 293 the cards to the response to one of the three spatial locations) given the top-down rule
 294 information from the PFC module? We sought to identify the structures in the input,
 295 recurrent and output connections of the sensorimotor module that give rise to this function.
 296

297 We start by observing that excitatory neurons in the sensorimotor module show a
 298 continuum of encoding strengths for task rule, response location and card features, and
 299 many neurons show conjunctive selectivity for these variables (Figure 5b, Supplementary
 300 Figure 6b). Therefore, we assigned each excitatory neuron in the sensorimotor module a
 301 preferred rule R , a preferred response location L and a preferred shared feature between

Figure 4: (Previous page.) **Structured top-down connections enable the propagation of the rule information.**

a. Each line represents the mean connection strength onto one excitatory neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent mean across neurons. PFC excitatory neurons project more strongly to excitatory neurons in the sensorimotor module that prefer the same rule (Student's t test, $p < .001$).

b. Each line represents the mean connection strength onto one VIP neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent mean across neurons. PFC excitatory neurons project more strongly to VIP neurons in the sensorimotor module that prefer the same rule (Student's t test, $p = .002$).

c. Each line represents the mean connection strength onto one SST neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent mean across neurons. PFC excitatory neurons project more strongly to SST neurons in the sensorimotor module that prefer the same rule (Student's t test, $p < .001$).

d. Each line represents the mean connection strength onto one excitatory neuron of the sensorimotor module, from the local SST neurons that prefer the same rule and the different rule. Bars represent mean across neurons. Local SST neurons project more strongly to excitatory neurons in the sensorimotor module that prefer the opposite rule (Student's t test, $p < .001$).

e. Each line represents the mean connection strength onto one PV neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent mean across neurons. PFC excitatory neurons project more strongly to PV neurons in the sensorimotor module that prefer the same rule (Student's t test, $p = 0.004$).

f. Each line represents the mean connection strength onto one excitatory neuron of the sensorimotor module, from the local PV neurons that prefer the same rule and the different rule. Bars represent mean across neurons. PV neurons in the sensorimotor module project more strongly to local excitatory neurons that prefer the opposite rule (Student's t test, $p < .001$).

g. The structure of the top-down connections as indicated by the results in **a-f**. The weaker connections are not shown.

Results in **a-f** are shown for an example network with subtractive dendritic nonlinearity. Networks with divisive and subtractive dendritic nonlinearity show similar patterns (Supplementary Figure 7).

302 the reference card and the test card at L , which we call F . For example, neurons with
303 $R = \text{color rule}$, $L = 1$ and $F = \text{blue}$ would have the highest activity during color rule
304 trials when the correct response is to choose the test card at location 1, and when that
305 test card shares the blue color with the reference card (it belongs to the population in the
306 sensorimotor module with the filled green color in Figure 5a). Intuitively, for this group of

307 neurons to show such selectivity, they should receive strong input from the input neurons
308 that encode the $F = blue$ feature of the test card at location $L = 1$ and the same feature
309 for the reference card. This would enable them to detect when the test card at $L = 1$ and
310 the reference card both have the feature $F = blue$.

311 In general, for neurons that prefer rule R , response location L and shared feature F ,
312 we can define their “preferred features” as the feature F of the reference card and the same
313 feature F for the test card at location L . Across all neurons, we found that the weights
314 from the input neurons encoding these preferred features were significantly stronger than
315 those encoding other features (Figure 5c). In addition, there is also an intuitive structure
316 in the output connections, where excitatory neurons in the sensorimotor module that prefer
317 a given response location L send stronger output connections to the output neuron that
318 prefers the same response location (Figure 5d). These structures were found to be consistent
319 across trained networks (Supplementary Figure 8).

320 These structures in the input and output connections give rise to an intuitive explana-
321 tion of how the sensorimotor module can perform rule-dependent action selection required
322 for the WCST. Here we illustrate this mechanism with an example trial (Figure 5a), where
323 the current rule is color, the reference card is a blue circle, and the test cards at locations 1
324 2 and 3 are blue triangle, red circle and red triangle, respectively. According to the rule
325 of WCST, the correct response should be location 1, since the test card at that location
326 matches the reference card in color. This choice can be generated as follows: the excitatory
327 population in the sensorimotor module that prefers $R = color\ rule$, $L = 1$ and $F = blue$ will
328 be most strongly activated since they not only receive strong top-down input from the PFC
329 module, but also the strongest feedforward input. Therefore, they are the most strongly
330 activated population (Figure 5a). And since they prefer response location 1, they will excite
331 the output neuron that prefers response location $L = 1$, which is the correct choice.

DISSECTING MODULAR RNNs TRAINED TO PERFORM A WCST ANALOG 17

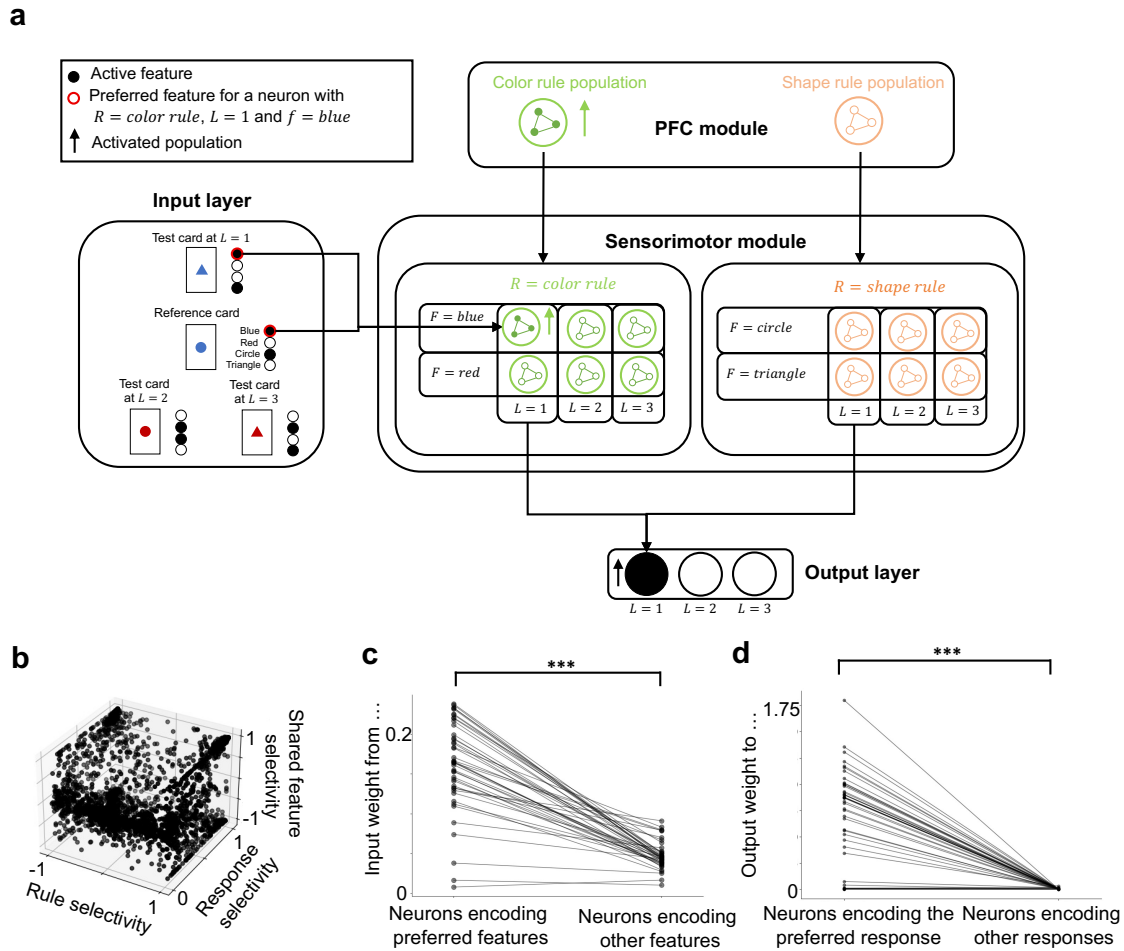


Figure 5: (Caption next page.)

Figure 5: (Previous page.) **Structures in the input and output weights of the sensorimotor module enable rule-dependent action selection.**

a. Excitatory neurons in the sensorimotor module were classified according to their preferred rule R , response location L and shared feature F . For example, neurons with $R = \textit{color rule}$, $L = 1$ and $F = \textit{blue}$ have the highest activity during color rule trials when the network chooses the test card at $L = 1$, and when that card shares the $F = \textit{blue}$ feature with the reference card. For a neuron with a given R , L and F , its “preferred features” are defined as the feature F of the reference card and same feature of the test card at location L . For example, the preferred features for neurons with $R = \textit{color rule}$, $L = 1$ and $F = \textit{blue}$ are the blue feature of the reference card and the test card at $L = 1$.

b. The joint distribution of the selectivity for rule (R), response location (L) and shared feature (F) across all neurons in the sensorimotor module. Result is aggregated across all trained networks.

c. Excitatory neurons in the sensorimotor module receive stronger connections from the input neurons that encode their preferred features (as defined in **a**). Student’s t-test, $p < .001$

d. Excitatory neurons in the sensorimotor module send stronger connections to the output neuron that represents their preferred response location. Student’s t-test, $p < .001$.

Panel **a** shows an example trial illustrating how the sensorimotor module can generate the correct response. During this trial, the reference card is a blue circle, and the test cards at location 1, 2, 3 are blue triangle, red circle and red triangle, respectively. The current rule is color. Therefore the correct response location should be $L = 1$. The network can generate that response because (a) the PFC population that encode the color rule are active, which send strong top-down excitation to the $R = \textit{color rule}$ population in the sensorimotor module; (b) the input neurons that encode the $F = \textit{blue}$ feature of the reference card and the test card at $L = 1$ are both active, which provide strong feedforward input to the excitatory population in the sensorimotor module with $R = \textit{color rule}$, $L = 1$ and $F = \textit{blue}$. Therefore, this population will be most strongly activated. Since they send strong long-range excitations to the output neuron that represents $L = 1$, the correct output neuron will be activated.

332 **Recurrent connectivity and dynamics within the sensorimotor module**

333 Given that different populations of neurons in the sensorimotor module receive dif-
334 ferential inputs about the external sensory stimuli and rule via the structured input and
335 top-down connections, how are they recurrently connected to produce dynamics that lead
336 to a categorical choice? To answer this, we first visualized the population neural dynamics
337 in the sensorimotor module by using principal component analysis (Figure 6a-b). As shown
338 in Figure 6a, neural trajectories during the inter-trial interval are clustered according to
339 the task rule. During the response period, the neural trajectories are separable according
340 to the response locations, albeit only in higher-order principal components (Figure 6b). In
341 addition, the subspaces spanned by neural trajectories of different rules and response loca-
342 tions are more orthogonal to each other compared to randomly shuffled data (Figure 6c-d,
343 Methods).

344 What connectivity structure gives rise to this signature in the population dynamics?
345 To answer this, we examined the pattern of connection weights between excitatory and
346 PV neurons that prefer different rules (R), response locations (L), and shared features (F)
347 by computing the connectivity biases between populations of neurons that are selective to
348 different rules (Figure 6e), response locations (Figure 6f) and shared features (Figure 6g). A
349 greater than zero connectivity bias means populations that prefer different rules (or response
350 locations or shared features) form a winner-take-all circuit structure analogous to the one
351 observed between rule-selective populations in the PFC module (c.f. top panel of Figure 3d,
352 details about how the connectivity biases were computed is described in Methods). We
353 observed that many of the connectivity biases were significantly above zero (Figure 6e-g),
354 especially for the ones that correspond to the inhibitory connections originating from the
355 PV neurons. This indicates that populations of neurons in the sensorimotor module that
356 are selective to different rules, response locations and shared features overall inhibit each
357 other. This mutual inhibition circuit structure magnifies the difference in the amount of
358 long-range inputs that different populations receive (Figure 5) and lead to a categorical
359 choice.

360 **SST neurons are essential to dendritic top-down gating**

361 The previous sections elucidate the key connectivity structures that enable the net-
362 work to perform the WCST. In this final section we are going to take advantage of the
363 biological realism of the trained RNN and examine the function of SST neurons in this
364 task.

365 It has been observed that different dendritic branches of the same neuron can be tuned
366 to different task variables [39, 40, 41, 42]. This property may enable individual dendritic
367 branches to control the flow of information into the local network [17, 20]. Given these

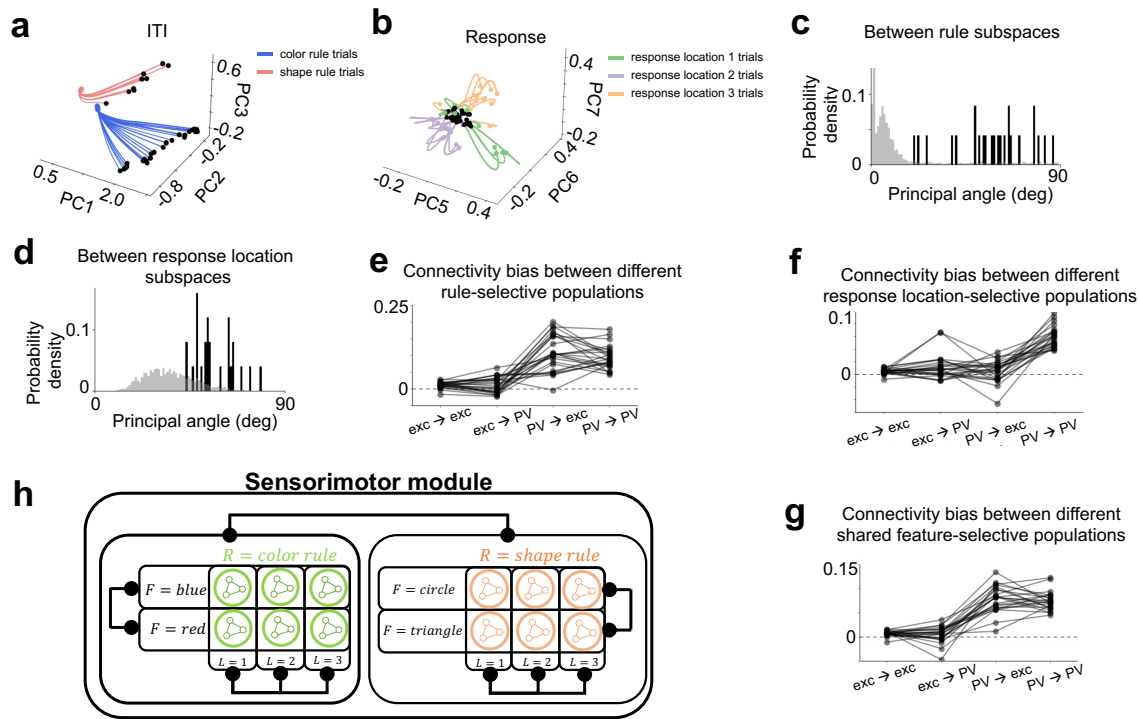


Figure 6: Recurrent dynamics and connectivity within the sensorimotor module.

a. Neural trajectories during the intertrial interval (ITI) for different task rules, visualized in the space spanned by the first three principal components. Black circles represent the start of the ITI. Only trials following a correct trial were included.

b. Neural trajectories during the response period for different choices, visualized in the space spanned by higher order principal components. Black circles represent the start of the response period. Only trials following a correct trial were included.

c. The principal angle between the subspaces spanned by neural trajectories during different task rules (gray distribution represents the principal angle obtained through shuffled data, see Methods). Each data point represents one trained network.

d. The principal angle between the subspaces spanned by neural trajectories during different responses (gray distribution represents the principal angle obtained through shuffled data, see Methods). Each data point represents one trained network.

e. The connectivity biases between different rule-selective populations across models.

f. The same as **e** but for different response location-selective populations.

g. The same as **e** but for different shared feature-selective populations.

h. The results in **e - g** show that neural populations selective for different rules, response locations and shared features mutually inhibit each other.

Data in **c-g** are shown for networks with subtractive dendritic nonlinearity. Networks with divisive dendritic nonlinearity show similar result (Supplementary Figure 9).

368 previous findings, we examined the coding of the top-down rule information at the level
369 of individual dendritic branches. Since each excitatory neuron in our networks is modeled
370 with two dendritic compartments, we examined the encoding of rule information by different
371 dendritic branches of the same excitatory neuron in the sensorimotor module.

372 One strategy of gating is for different dendritic branches of the same neuron to prefer
373 the same rule, in which case these neurons form distinct populations that are preferentially
374 recruited under different task rules (population-level gating, Figure 7a, right). An alterna-
375 tive strategy is for different dendritic branches of the same neuron to prefer different rules,
376 which would enable these neurons to be involved in both task rules (dendritic branch-specific
377 gating, Figure 7a, left).

378 In light of this, we examined for our trained networks to what extent they adopt these
379 strategies. We found that the rule selectivity between different dendritic branches of the
380 same neuron were highly correlated (Figure 7b). This indicates that the trained networks
381 are mostly using the population-level gating strategy, where different dendritic branches of
382 the same neuron encode the same rule.

383 What factors might determine the extent to which the trained networks adopt these
384 two strategies? Previous modeling work suggests that sparse connectivity from SST neurons
385 to the dendrites of the excitatory neurons increases the degree of dendritic branch-specific
386 gating, in the case where the connectivity is random (Figure 4f in Ref.[20]). To see if the
387 same effect is present in our task-optimized network with structured connectivity, we re-
388 trained networks with different levels of sparsity from 0 to 0.8 and studied its effect on
389 the dendritic branch specificity of rule coding (Methods). We found that the degree of
390 dendritic branch-specific encoding of the task rule increased with sparsity (see Figure 7c,
391 d for subtractive dendritic nonlinearity; Supplementary Figure 10a for divisive dendritic
392 nonlinearity). Intuitively, when the connection is sparse, a smaller number of SST neurons
393 target each dendritic branch, making it more likely that the branch receives an uneven
394 number of inputs from SST neurons selective for different rules. Taken together, we observed
395 that the trained networks adopted a mixture of population-level and dendritic-level gating
396 strategies for top-down control, and the balance between the two strategies depends on the
397 sparsity of the connections from the SST neurons to the dendrites of excitatory neurons.

398 Indeed, SST neurons play a causal role in relaying the top-down rule information into
399 the sensorimotor network and reconfiguring its dynamics according to the task rule. We
400 simulated optogenetic inhibition by silencing the SST neurons in the sensorimotor module,
401 which significantly impaired task performance (Figure 7e, see Methods section for details of
402 the protocol). In addition, the principal angle between the subspaces for different rules (Fig-
403 ure 6c) significantly decreased after SST neurons in the sensorimotor module were silenced
404 (Figure 7f). Silencing of the SST neurons in the sensorimotor module also significantly

405 diminished nonlinear mixed-selective coding of rule and stimulus among the excitatory neu-
406 rons in the sensorimotor module (Figure 7g, Supplementary Figure 11, Methods), which has
407 been proposed to be important for rule-based sensorimotor associations [43, 44, 45]. Taken
408 together, these results highlight the role that SST neurons in the sensorimotor module play
409 during top-down control. This analysis also shows that by combining artificial neural net-
410 work with knowledge from neurobiology, it is possible to probe the functions of fine-scale
411 biological components in cognitive behaviors.

412 Discussion

413 In this paper, we analyzed recurrent neural networks trained to perform a classic
414 task involving un-cued task switching - the Wisconsin Card Sorting Test. The networks
415 consist of a “PFC” module trained to represent the rule and interacts with a “sensorimotor”
416 module that instantiates different sensorimotor mappings depending on the rule. In order
417 to study the functions of dendritic computation and different neuronal types, each module
418 is endowed with excitatory neurons with two dendritic branches as well as three major types
419 of inhibitory neurons - PV, SST and VIP. After training, we dissected the trained networks
420 to elucidate a number of intra-areal and inter-areal neural circuit mechanisms underlying
421 WCST, as summarized in Figure 8.

422 Mapping between model components and brain regions

423 Different components of the trained network can be mapped to different brain regions
424 (Figure 8). While single neurons in the dorsal-lateral PFC (DLPFC) are shown to encode the
425 task rule [46], neurons in the anterior cingulate cortex (ACC) are thought to be important
426 for performance monitoring [47], and have been shown to receive more input about the
427 feedback [48, 49, 50, 51]. Therefore, the rule neurons and conjunctive error x rule neurons
428 in the model correspond to the putative functions of the neurons in DLPFC and ACC. The
429 input to the PFC module about negative feedback may come from subcortical areas such
430 as the amygdala [52] or from the dopamine neurons in the substantia nigra pars compacta
431 (SNc) and ventral tegmental area (VTA) [53, 54]. The sensorimotor module may correspond
432 to parietal cortex or basal ganglia which have been shown to be involved in sensorimotor
433 transformations [55, 56]. The neurons in the input layer that encode the color and shape of
434 the card stimuli exist in higher visual areas such as the inferotemporal cortex [57, 58, 59].
435 The neurons in the output layer that encode different response locations could correspond
436 to movement location-specific neurons in the motor cortex [60].

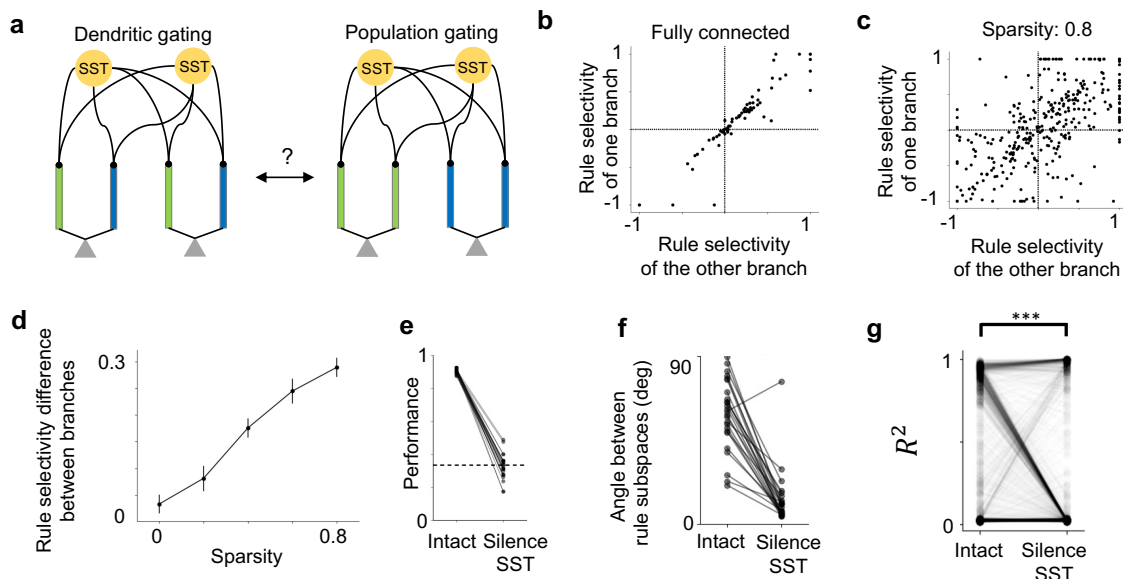


Figure 7: Examining the role of SST neurons in the sensorimotor module in top-down gating.

a. Two scenarios for top-down gating. Blue and green color represent dendritic branches that prefer one of the two rules. Different dendritic branches of the same neuron could have similar (right) or different (left) rule selectivity.

b. The rule selectivity of one dendritic branch against the other, aggregated across all models where the connections from the SST neurons to the excitatory neurons are all-to-all. The rule selectivity for different dendritic branches of the same neuron are highly correlated.

c. The rule selectivity of one dendritic branch against the other, aggregated across all models where the 80% of the connections from the SST neurons to the excitatory neurons are frozen at 0 throughout training. Note the the rule selectivity for different dendritic branches of the same neuron are less correlated than in **b**.

d. The degree of dendritic branch-specific encoding of the task rule is quantified as the difference in the rule selectivity between the two dendritic branches of the same excitatory neuron in the sensorimotor module. Across all dendritic branches, this quantity increases with the sparsity of the SST \rightarrow dendrite connectivity.

e. Task performance drops significantly after silencing SST neurons in the sensorimotor module. Each line represents a trained network.

f. The principal angle between rule subspaces (c.f. Figure 6c) drops significantly after silencing SST neurons in the sensorimotor module. Each line represents a trained network.

g. The strength of conjunctive coding of rule and stimulus (as measured by the R^2 value in a linear model with conjunctive terms, see Methods) decreases after silencing SST neurons in the sensorimotor module (Student's t-test, $p < .001$). Each line represents one neuron. Results are aggregated across networks.

Results in **b-g** are for networks with subtractive dendritic nonlinearity. Networks with divisive dendritic nonlinearity show similar result (Supplementary Figure 10).

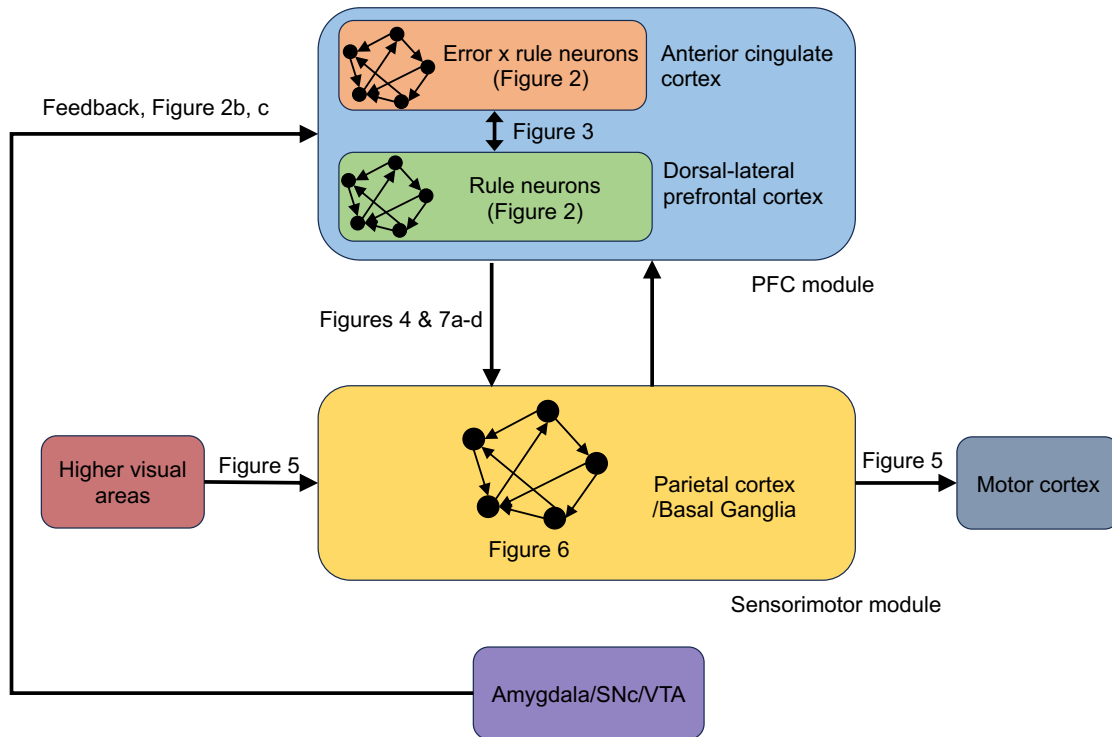


Figure 8: A summary of the main results.

Different components of the model can be mapped to different brain regions; The conjunctive error x rule neurons may reside in the anterior cingulate cortex; The rule neurons may be found in the dorsal-lateral PFC; The input to the PFC module about negative feedback may come from subcortical areas such as the amygdala or the midbrain dopamine neurons; The sensorimotor module may correspond to parietal cortex or basal ganglia which have been shown to be involved in sensorimotor transformations; Neurons in the input layer that encode the color and shape of the card stimuli exist in higher visual areas such as the inferotemporal cortex; Neurons in the output layer that encode different response locations could correspond to neurons in the motor cortex.

437 **Attractor states supported by inter-areal connections.**

438 We observed that in many networks, the interaction between the two modules was
439 needed to sustain the two rule attractor states (Supplementary Figure 2b,c), although the
440 majority of the excitatory input to PFC neurons come from local population (Supplementary
441 Figure 2d). Traditionally, it was thought that local interactions within the frontal cortex are
442 sufficient for the maintenance of the persistent activity [28, 29, 30, 31]. Recent large-scale
443 electrophysiological recordings, on the other hand, revealed highly distributed encoding of
444 cognitive variables [32, 34, 61, 62, 63, 64]. In addition, distributed patterns of persistent
445 activity emerge in neural network models of multiple brain regions that are constrained
446 by anatomical and neurophysiological data [65, 66]. Despite the empirical evidence, the

447 functional advantages of this multi-areal encoding scheme remain an open question.

448 **Circuit mechanism in the frontal-parietal network for rule maintenance and** 449 **update**

450 Two distinct types of responses among the excitatory neurons emerge in the PFC
451 module as a result of training: neurons that only encode the rule, and neurons that con-
452 junctively encode negative feedback and rule. Neurons that show conjunctive selectivity
453 for rule and negative feedback have been reported in monkey prefrontal and parietal cor-
454 tices while they perform the same WCST task [8, 9]. Theoretical work suggests that these
455 mixed-selective neurons are essential if the network needs to switch between *different* rule
456 attractor states after receiving the *same* input that signals negative feedback [67].

457 We further revealed the connectivity pattern between different populations of excita-
458 tory and PV neurons in the PFC module in order for the network to switch between rule
459 attractor states (Figure 3c). In addition, this connectivity pattern is consistent across dozens
460 of trained networks with different initializations and dendritic nonlinearities (Figure 3d and
461 Supplementary Figure 4). This circuit mechanism bears resemblance to a previous circuit
462 model of WCST [68]. In that model, the switching between different rule states is achieved
463 by synaptic desensitization caused by the convergence of two signals - one that signals the
464 recent activation of the synapse, and another that signals the negative feedback. However,
465 that model does not predict the existence of neurons with conjunctive coding of negative
466 feedback and rule, which has been observed experimentally [8, 9].

467 The simplified circuit for the PFC module in Figure 3c can be applied not only to
468 rule switching, but to the switching between other behavioral states as well. For example, it
469 resembles the head-direction circuit in fruit fly [69], where the offset in the connections be-
470 tween the neurons coding for head direction and those coding for the conjunction of angular
471 velocity and head direction enables the update of the head-direction attractor state by the
472 angular velocity input. In addition, this circuit structure may underlie the transition from
473 staying to switching during patch foraging behavior. Indeed, in a laboratory task mimicking
474 natural foraging for monkeys, it was found that neurons in the anterior cingulate cortex
475 increase their firing rates to a threshold before animals switch to another food resource [70].

476 **Connecting subspace to circuits**

477 Methods that describe the representation and dynamics on the neuronal population
478 level have gained increasing popularity and generated novel insights that cannot be discov-
479 ered using single neuron analysis (e.g. [60, 71]). In the meantime, it would be valuable to
480 connect population-level phenomena to their underlying circuit basis [72]. In our model, we
481 found that silencing of the SST neurons has a specific effect on the population-level repre-

482 sentation, namely, it decreases the angle between rule subspaces (Figure 7f). We also found
483 that silencing the other types of inhibitory neurons has different effects (data not shown).
484 Silencing the PV neurons to an instability of the network dynamics, whereas silencing the
485 VIP neurons caused an insignificant decrease of the network performance. The lack of ef-
486 fect after silencing the VIP neurons is due to the fact that the VIP neurons were largely
487 inhibited by the SST neurons in the trained model. Future work could study the function
488 of VIP neurons under different connectivity constraints between SST and VIP neurons.

489 **Dynamics of behavior during rule switching**

490 Our networks switch rules in just one trial (Figure 1c, d). This fast switching agrees
491 with the monkey behavior in some studies [21, 9, 27], but other studies report that monkeys
492 switch rules using on average tens of trials [8, 6]. For example, in Ref.[6], monkeys' perfor-
493 mance is at chance after a single error (Figure 3D in [6]), and they gradually use positive
494 feedback to reinforce their behavior according to the new rule (Figure 4A in [6]). However,
495 when our network model was trained to achieve less than perfect accuracy, switching after
496 a rule change now takes a few trials (Figure 1f) similar to behavioral observations of many
497 monkey experiments. In this case, the dissected network mechanism is compatible with
498 what we reported here, albeit not as clear cut.

499 Indeed, in WCST and related rule switching paradigms, subjects' performance is
500 often not perfect even during trials when the rule is fixed. This is possibly because during
501 training, the rule is switched when the subjects' performance reaches a certain criterion
502 (e.g. 85% correct in a sequence of 20 trials in Ref.[6]). In that case, negative feedback can
503 due to either a rule switch or the inaccuracy in the sensorimotor transformation (even under
504 the correct rule). Therefore, subjects need to integrate information across several trials to
505 decide whether the rule has actually switched. For example, Purcell and Kiani [73] analyzed
506 the behavior of humans in an environment switching task. The task is similar to the WCST
507 analog used in this paper, except that the noise level in the stimuli varies from trial to
508 trial. It was shown that the behavior of subjects can be well described by a Bayesian ideal
509 observer model, where the evidence towards an environment switch is increased whenever
510 the subjects performs an error trial, and the amount of increase depends on the difficulty
511 of that trial: the easier the error trial is, the more likely that the environment has switched
512 and the larger the incremental evidence towards an environment switch is.

513 Aside from the difficulty of the task under a fixed rule, the relationship between the
514 different rules may also play a role in how fast animals can switch between them. In tasks
515 that involve simple reversal of motor response or sensorimotor mappings, monkeys usually
516 use a small number of trials to switch between rules [74, 75, 76]. On the other hand, for
517 WCST with more than two rules, as is usually used for humans, subjects typically use more

518 trials to switch to the new rule [1, 77].

519 There are other reasons that may contribute to the suboptimality of behavior during
520 rule switching, including random exploration [77], poor sensitivity to negative feedback [77],
521 integration of reward history across multiple trials [78, 73, 10, 79], the gradual update of
522 the value of the counterfactual rule [80] or the cost of cognitive control [81]. Neuronal
523 mechanisms on longer timescales such as synaptic mechanisms [82] may be required to
524 produce the slow switching behavior.

525 In conclusion, our approach of incorporating neurobiological knowledge into train-
526 ing RNNs can provide a fruitful way to build circuit models that are functional, high-
527 dimensional, and reflect the heterogeneity of biological neural networks. In addition, dis-
528 secting these networks can make useful cross-level predictions that connect biological ingre-
529 dients with circuit mechanisms and cognitive functions.

530 **Acknowledgements:** This work was supported by James Simons Foundation Grant
531 543057SPI, the National Institutes of Health grant R01MH062349, and the ONR grant
532 N00014-23-1-2040. YL thanks (alphabetically) Aldo Battista, Mark Buckley, Sage Chen,
533 Shuo Chen, Vishwa Goudar, Kenneth Kay, Haohong Li's lab, Jianguang Ni's lab, Yu Qi's
534 lab, Yi Sun's lab, Lucas Tian, Bo Shen, Xiaohan Zhang and all members of Xiao-Jing
535 Wang's lab for helpful discussions and comments on the manuscript.

536 **Data availability statement:** All computer code used to generate the model and results
537 will be uploaded to Github upon acceptance of the manuscript.

538 Methods

539 Model setup

540 The RNN consists of two bidirectionally-connected modules, the PFC module and
541 the sensorimotor module. Each module consists of 70 excitatory neurons and 30 inhibitory
542 neurons. Each excitatory neuron has 2 dendritic compartments. The inhibitory neurons
543 are evenly divided into three types: PV, SST and VIP. Different types of neurons have
544 different connectivity, inspired by experimental findings [83]: PV neurons target the somatic
545 compartment of excitatory neurons and other PV neurons, SST neurons target the dendritic
546 compartment of excitatory neurons as well as PV and VIP neurons, and VIP neurons target
547 SST neurons. Excitatory neurons target other excitatory neurons, PV and SST neurons.
548 The connection strength between all other types of neurons were fixed at zero throughout
549 training.

550 Only excitatory neurons send long-range projections to other modules. The long-
551 range projections from the sensorimotor module to the PFC module target the dendritic
552 compartment of the excitatory neurons and the PV neurons. This is inspired by the ex-
553 perimental evidence that PV neurons mediate feedforward inhibition [12]. The long-range
554 top-down projections from the PFC to the sensorimotor module target the dendritic com-
555 partments of the excitatory neurons and all three types of inhibitory neurons. Finally,
556 external inputs to both modules target the dendritic compartment of excitatory neurons
557 and PV neurons.

558 The dynamics of the somata of the excitatory neurons in the RNN are described by

$$\tau \frac{dh_{\text{esoma}}}{dt} = -h_{\text{esoma}} + f_{\text{soma}}(W_{\text{esoma} \rightarrow \text{esoma}}^{\text{rec}} h_{\text{esoma}} + W_{\text{PV} \rightarrow \text{esoma}}^{\text{rec}} h_{\text{PV}} + \sum_{\text{dendrites}} h_{\text{dendrite}}), \quad (1)$$

559 where $\tau = 100$ ms, $dt = 10$ ms. Somata of excitatory neurons in both the sensorimotor
560 and PFC modules obey the same equation. Here “esoma” stands for the soma of excita-
561 tory neurons. $W_{\text{esoma} \rightarrow \text{esoma}}^{\text{rec}}$ and $W_{\text{PV} \rightarrow \text{esoma}}^{\text{rec}}$ represent the connectivity matrix between
562 the soma of excitatory neurons and from the local PV neurons to the soma of excitatory
563 neurons, respectively. h_{esoma} and h_{PV} are the activity of the soma of excitatory neurons
564 and PV neurons. h_{dendrite} is the activity of the dendritic compartment. f_{soma} is the somatic
565 nonlinear activation function which was modeled as a rectified linear function:
566

$$f_{\text{soma}} = \begin{cases} x, & x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

568 The dendritic activity is a nonlinear function of the excitatory and inhibitory inputs.

$$h_{\text{dendrite}} = f_{\text{dendrite}}(I_{\text{exc}}, I_{\text{inh}}). \quad (3)$$

570 I_{exc} is the total excitatory input to the dendrite. It consists of long-range inputs from the
571 input neurons (neurons that encode the feedback for the PFC module and neurons that
572 encode the stimulus for the sensorimotor module) as well as the long-range input from the
573 excitatory neurons in the other module. $I_{\text{exc}} = I_{\text{in}} + I_{\text{cross-module}}$. I_{inh} is the inhibitory input
574 to the dendrite from the local SST neurons. $I_{\text{inh}} = I_{\text{SST} \rightarrow \text{edend}}$. Here “edend” stands for
575 the dendrite of excitatory neurons. The functional form of f_{dendrite} is described in the next
576 section.

577 The inhibitory neurons are modeled as standard point neurons. Different types of
578 inhibitory neurons receive different input connections. In the sensorimotor (SM) module,

579 the dynamics of PV neurons are described by

$$\begin{aligned}
 \tau \frac{dh_{SM,PV}}{dt} = & -h_{SM,PV} + f_{soma}(W_{SM,PV \rightarrow SM,PV}^{rec} h_{SM,PV} \\
 & + W_{SM,SST \rightarrow SM,PV}^{rec} h_{SM,SST} \\
 & + W_{SM,esoma \rightarrow SM,PV}^{rec} h_{SM,esoma} \\
 & + W_{in \rightarrow SM,PV} u_{sensory} \\
 & + W_{PFC,esoma \rightarrow SM,PV} h_{PFC,esoma}),
 \end{aligned} \tag{4}$$

581 where $W_{SM,PV \rightarrow SM,PV}^{rec}$, $W_{SM,SST \rightarrow SM,PV}^{rec}$, $W_{SM,esoma \rightarrow SM,PV}^{rec}$ are the connection weight
582 matrices between the PV neurons, from local SST neurons to the PV neurons, and from local
583 excitatory neurons to the PV neurons, respectively. $W_{in \rightarrow SM,PV}$ is the input weight matrix
584 to the PV neurons, and $u_{sensory}$ is the input to the sensorimotor module that represents
585 the features about the cards. $W_{PFC,esoma \rightarrow SM,PV} h_{PFC,esoma}$ is the top-down connection
586 weight matrix from the excitatory neurons in the PFC module to the PV neurons in the
587 sensorimotor module.

588 For the SST neurons,

$$\begin{aligned}
 \tau \frac{dh_{SM,SST}}{dt} = & -h_{SM,SST} + f_{soma}(W_{SM,VIP \rightarrow SM,SST}^{rec} h_{SM,VIP} \\
 & + W_{SM,esoma \rightarrow SM,SST}^{rec} h_{SM,esoma} \\
 & + W_{PFC,esoma \rightarrow SM,SST} h_{PFC,esoma}),
 \end{aligned} \tag{5}$$

590 where $W_{SM,VIP \rightarrow SM,SST}^{rec}$ and $W_{SM,esoma \rightarrow SM,SST}^{rec}$ are the connection weight matri-
591 ces from local VIP neurons and excitatory neurons to the SST neurons, and
592 $W_{PFC,esoma \rightarrow SM,SST} h_{PFC,esoma}$ is the top-down connection weight matrix from the exci-
593 tatory neurons in the PFC module to the SST neurons in the sensorimotor module.

594 For the VIP neurons,

$$\begin{aligned}
 \tau \frac{dh_{SM,VIP}}{dt} = & -h_{SM,VIP} + f_{SM,soma}(W_{SM,SST \rightarrow SM,VIP}^{rec} h_{SM,SST} \\
 & + W_{PFC,esoma \rightarrow SM,VIP} h_{PFC,esoma}),
 \end{aligned} \tag{6}$$

596 where $W_{SM,SST \rightarrow SM,VIP}^{rec}$ is the connection weight matrix from the local SST neurons to the
597 VIP neurons, and $W_{PFC,esoma \rightarrow SM,VIP} h_{PFC,esoma}$ is the top-down connection weight matrix
598 from the excitatory neurons in the PFC module to the VIP neurons in the sensorimotor
599 module.

600 The inhibitory neurons in the PFC module are described by similar equations, except
601 only the PV neurons receive long-range bottom-up inputs from the sensorimotor module:

$$\begin{aligned} \tau \frac{dh_{\text{PFC,PV}}}{dt} = & -h_{\text{PFC,PV}} + f_{\text{soma}}(W_{\text{PFC,PV} \rightarrow \text{PFC,PV}}^{\text{rec}} h_{\text{PFC,PV}} \\ & + W_{\text{PFC,SST} \rightarrow \text{PFC,PV}}^{\text{rec}} h_{\text{PFC,SST}} \\ & + W_{\text{PFC,esoma} \rightarrow \text{PFC,PV}}^{\text{rec}} h_{\text{PFC,esoma}} \\ & + W_{\text{in} \rightarrow \text{PFC,PV}} u_{\text{feedback}} \\ & + W_{\text{SM,esoma} \rightarrow \text{PFC,PV}} h_{\text{SM,esoma}}), \end{aligned} \quad (7)$$

$$\begin{aligned} \tau \frac{dh_{\text{PFC,SST}}}{dt} = & -h_{\text{PFC,SST}} + f_{\text{soma}}(W_{\text{PFC,VIP} \rightarrow \text{PFC,SST}}^{\text{rec}} h_{\text{PFC}} \\ & + W_{\text{PFC,esoma} \rightarrow \text{PFC,SST}}^{\text{rec}} h_{\text{PFC,esoma}}), \end{aligned} \quad (8)$$

$$\tau \frac{dh_{\text{PFC,VIP}}}{dt} = -h_{\text{PFC,VIP}} + f_{\text{soma}}(W_{\text{PFC,SST} \rightarrow \text{PFC,VIP}}^{\text{rec}} h_{\text{PFC,SST}}), \quad (9)$$

where u_{feedback} represents the external input to the PFC module about the feedback of the previous trial.

In practice, we used a mask matrix to enforce the connectivity between different cell types.

$$W^{\text{rec}} = |\tilde{W}^{\text{rec}}| * M + W^{\text{fix}}, \quad (10)$$

where \tilde{W}^{rec} is the unconstrained connectivity matrix updated by the learning algorithm, M is a matrix consisting of 1, 0 and -1 depending on whether the corresponding connection is excitatory, inhibitory or nonexistent. W^{fix} implements the fixed coupling between the dendrite and the soma.

Only the somata of excitatory neurons send output connections. The output in each module is generated via a simple linear readout:

$$y_{\text{SM}} = W_{\text{out, SM}} h_{\text{SM, esoma}}. \quad (11)$$

$$y_{\text{PFC}} = W_{\text{out, PFC}} h_{\text{PFC, esoma}}. \quad (12)$$

The output connections were constrained to be positive.

Variations in the model hyperparameters

Dendritic nonlinearities. We trained models with two types of dendritic nonlinearities f_{dendrite} - subtractive and divisive. They are inspired by in-vitro and computational studies showing different types of inhibitory modulation on the dendritic activity depending on the location of inhibition relative to excitation [23]. Both types of dendritic nonlinearities are sigmoidal functions of the excitatory input. Under subtractive nonlinearity, as the

625 inhibitory input increases, the turning point of the sigmoid function moves to larger values,
626 consistent with the experimental observation when the inhibitory current is injected at the
627 same location or more distal than the excitation [23]. For the divisive nonlinearity, the turn-
628 ing point of the sigmoid is not affected by the level of inhibition, but the saturating level of
629 the sigmoid function decreases with the level of inhibition, consistent with the experimental
630 observation when the inhibitory current is injected close to the soma [23].

631 The equations of the different dendritic nonlinearities are given by:

$$632 \quad f_{\text{dendrite}}^{\text{subtractive}}(I_{\text{exc}}, I_{\text{inh}}) = \tanh(I_{\text{exc}} - I_{\text{inh}})$$

$$633 \quad f_{\text{dendrite}}^{\text{divisive}}(I_{\text{exc}}, I_{\text{inh}}) = k_1(1 + \tanh(I_{\text{exc}} - 1)) + k_2,$$

634 where $k_1 = \frac{1}{e^{I_{\text{inh}}}}$ and $k_2 = -1 - \tanh(-1)$. The form of the divisive dendritic nonlinearity
635 was specified such that it is divisively modulated by I_{inh} (even when $I_{\text{exc}}=0$), and that it is
636 0 only when both I_{exc} and I_{inh} are 0.

637 **Initializations.** The connectivity matrices were initialized either using a normal
638 distribution with mean 0 and standard deviation $\sqrt{\frac{2}{N}}$ (where N is the total number of
639 recurrent units) or a uniform distribution between $-\sqrt{\frac{6}{N}}$ and $\sqrt{\frac{6}{N}}$.

640 **Sparsity of the SST→dendrite connectivity in the sensorimotor module.**

641 To study how the degree of dendritic branch-specific rule encoding in the sensorimotor
642 module is affected by the sparsity of the connections from SST neurons to the dendrite
643 of excitatory neurons, we varied this sparsity by fixing a fraction of randomly chosen con-
644 nections to be 0 throughout training. The sparsity levels used were 0, 0.2, 0.4, 0.6 and
645 0.8.

646 **Random seeds.** For each combination of the hyperparameter configuration intro-
647 duced above (except the sparsity), we trained models using 50 random seeds for Pytorch
648 (other random seeds were fixed). For each sparsity level other than 0, we trained models
649 using 10 random seeds for Pytorch.

650 Task

651 The networks were trained on an analog of the Wisconsin Card Sorting Test (WCST)
652 used for monkeys [21, 8, 9]. Each trial starts with the presentation of a “reference card” for
653 500 ms, after which three “test cards” appear around the reference card for 500 ms. Each
654 card contains an object with a specific color (blue or red) and shape (circle or triangle).
655 Among the three test cards, one of them matches the color of the reference card, another
656 one matches the shape of the reference card, and the third card matches neither feature
657 of the reference card. Depending on the rule (color or shape), the location where the test

658 card has the same color or shape feature as the reference card should be chosen. The choice
659 should be made during the 500 ms when both the reference card and the test cards are
660 presented. At the end of this period, a feedback signal is presented for 100 ms, indicating
661 whether the choice is correct or incorrect. This is followed by a 1 second inter-trial interval.

662 The task rule switches after a random number of trials, without informing the network.
663 Therefore, the network inevitably makes an error for the first trial after the rule switch since
664 it has not yet received the information that the rule has switched. The network should then
665 adjust its behavior to the new rule by utilizing the feedback signal.

666 Representation of inputs and outputs

667 Each card is represented as a four-dimensional binary vector, where different entries
668 represent the presence of the two colors and shapes. The feedback input is a two-dimensional
669 one-hot vector, where the two entries represent positive and negative feedback. The target
670 output for the sensorimotor module is a three-dimensional one-hot vector, where each entry
671 represents one response location on the screen. This target is non-zero only during the 500
672 ms response period when both the reference card and the test cards are presented. The
673 target output for the PFC module is a two-dimensional one-hot vector, where each entry
674 represents one rule. This target is non-zero during the entire trial.

675 Training method

676 During training, the networks ran continuously across 20 consecutive trials with 3
677 random rule switches. Importantly, the network dynamics were not reset during the inter-
678 trial interval. The loss function was aggregated across the 20 trials.

$$L = \sum_{trial=1}^{20} \sum_t \left(\|y_{PFC}(trial, t) - \hat{y}_{rule}(trial, t)\|^2 + \|y_{SM}(trial, t) - \hat{y}_{choice}(trial, t)\|^2 \right), \quad (13)$$

679
680 where $y_{PFC}(trial, t)$ and $y_{SM}(trial, t)$ are the activity of the readout neurons for the PFC
681 and sensorimotor module at time t in a given trial, respectively. $\hat{y}_{rule}(trial, t)$ is the target
682 output for the PFC module which represents the rule of the current trial. It is a binary
683 vector of dimension 2 where each entry represents one rule. The activation of the entry
684 that represents the correct rule is 1 throughout the entire trial. $\hat{y}_{choice}(trial, t)$ is the target
685 output for the sensorimotor module which represents the correct choice for the current trial.
686 It is a binary vector of dimension 3 where each entry represents one of the three response
687 locations. The activation of the entry that represents the correct choice is 1 during the
688 response period (500 ms when both the reference card and the tests card are shown).

689 The standard backpropagation through time algorithm [84] with the Adam optimizer
690 [85] was used to update all connection weights.

691 We also used curriculum learning to speed up training. Initially, the stimulus, choice
692 and outcome of the previous trial were all provided to the PFC module as input. This
693 way all the information needed to perform the current trial is contained within the input,
694 and the networks do not need to memorize past trials. During the training phase, the
695 network performs 20 consecutive trials with 3 random rule switches, therefore the maximum
696 performance is 85%. When the training performance reached above 65%, we started testing
697 the network on longer trial sequences (200 consecutive trials with 10 rule switches). The
698 maximum performance during testing is 95%. If the networks reached on average 90%
699 performance during the recent 5 tests, the input about the previous stimulus was removed.
700 When the networks reached 90% performance again, the information about the previous
701 choice information was removed. The networks were then trained until they reached 90%
702 performance.

703 Lower performance criteria was used for the model trained using early stopping (Fig-
704 ure 1f). In particular, curriculum training advanced to the next stage when the testing
705 performance reached 80%.

706 **Single neuron selectivity metric**

707 The selectivity index (SI) for rule is defined as

$$708 \quad SI_{\text{rule}} = \frac{h(\text{color}) - h(\text{shape})}{|h(\text{color})| + |h(\text{shape})|}, \quad (14)$$

709 where $h(\text{color})$ and $h(\text{shape})$ represent the trial-averaged single neuron activity during color
710 rule and shape rule, respectively. Neural activity was first averaged over the inter-trial
711 interval before further being averaged across trials.

712 The error selectivity is defined similarly

$$713 \quad SI_{\text{error}} = \frac{h(\text{after error}) - h(\text{after correct})}{|h(\text{after error})| + |h(\text{after correct})|}, \quad (15)$$

714 where $h(\text{after error})$ and $h(\text{after correct})$ are the mean single neuron activity after error and
715 correct trials, respectively. Neural activity was first averaged across the feedback presenta-
716 tion and inter-trial interval periods before being averaged across trials.

717 The selectivity for response location is defined as

$$718 \quad SI_{\text{response}} = \frac{h(L^*) - h(\bar{L})}{|h(L^*)| + |h(\bar{L})|}, \quad (16)$$

719 where L^* the preferred response location of the neuron, and $h(L^*)$ represents the mean
720 activity across trials when the network chooses location L^* . $h(\bar{L})$ represents the mean
721 activity across trials when the choice of the network is not location L^* . Therefore, this
722 selectivity index ranges from 0 to 1. We included neural activity during the response period
723 when computing this selectivity index.

724 Neurons that prefer color/shape rule were further divided according to their preferred
725 shared feature. The selectivity for the shared feature is defined as

$$726 \quad SI_{\text{shared feature}} = \frac{h(\text{blue}) - h(\text{red})}{|h(\text{blue})| + |h(\text{red})|}, \quad (17)$$

727 for neurons that prefer the color rule, and

$$728 \quad SI_{\text{shared feature}} = \frac{h(\text{circle}) - h(\text{triangle})}{|h(\text{circle})| + |h(\text{triangle})|}, \quad (18)$$

729 for neurons that prefer the shape rule. Here $h(\text{blue})$, $h(\text{red})$, $h(\text{circle})$, $h(\text{triangle})$ represent
730 the mean activity of a neuron across trials when the reference card is blue, red (when the
731 current rule is color), circle or triangle (when the current rule is shape). We included neural
732 activity during the response period when computing this selectivity index.

733 **Classification criteria for different neuronal populations**

734 Each neuron in the PFC module was classified as a “rule neuron” if the absolute value
735 of its rule selectivity was greater than 0.5 and the absolute value of its error selectivity was
736 smaller than 0.5. The rest of the neurons were classified as “error neurons” if their error
737 selectivity was greater than 0.5. Error neurons with greater mean activity during the color
738 rule trials that follow an error trial were defined as error x color rule neurons, and the other
739 error neurons were defined as error x shape rule neurons.

740 Each neuron in the sensorimotor module was assigned with a preferred rule, response
741 location and shared feature according to the condition during which it has the highest
742 activity. There was no threshold for this classification.

743 **Connectivity bias**

744 The connectivity bias (CB) was defined as the difference in the average connection
745 weight between different sub-population of neurons. A positive value indicates an agreement
746 with the simplified circuit diagram (Figures 3c, Figure 6h). For example, the connectivity
747 bias from the PFC PV neurons to the PFC excitatory neurons is given by

$$\begin{aligned}
 CB(\text{PFC PV} \longrightarrow \text{PFC E}) &= \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC E rule2}) \\
 &+ \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC E rule1}) \\
 &- \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC E rule1}) \\
 &- \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC E rule2}),
 \end{aligned} \tag{19}$$

749 where for example $\bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC E rule2})$ represents the average (unsigned)
 750 connection strength from the PFC PV neurons that prefer rule 1 to PFC excitatory neurons
 751 that prefer rule 2. Here rule 1 refers to color rule and rule 2 refers to shape rule.

752 The other connectivity biases were defined analogously.

$$\begin{aligned}
 CB(\text{PFC E} \longrightarrow \text{PFC E}) &= \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC E rule1}) \\
 &+ \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC E rule2}) \\
 &- \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC E rule2}) \\
 &- \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC E rule1})
 \end{aligned} \tag{20}$$

$$\begin{aligned}
 CB(\text{PFC E} \longrightarrow \text{PFC PV}) &= \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC PV rule1}) \\
 &+ \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC PV rule2}) \\
 &- \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC PV rule2}) \\
 &- \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC PV rule1})
 \end{aligned} \tag{21}$$

$$\begin{aligned}
 CB(\text{PFC PV} \longrightarrow \text{PFC PV}) &= \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC PV rule2}) \\
 &+ \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC PV rule1}) \\
 &- \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC PV rule1}) \\
 &- \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC PV rule2})
 \end{aligned} \tag{22}$$

$$\begin{aligned}
 CB(\text{PFC E} \longrightarrow \text{PFC error x rule}) &= \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC error x rule2}) \\
 &+ \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC error x rule1}) \\
 &- \bar{W}(\text{PFC E rule1} \longrightarrow \text{PFC error x rule1}) \\
 &- \bar{W}(\text{PFC E rule2} \longrightarrow \text{PFC error x rule2})
 \end{aligned} \tag{23}$$

$$\begin{aligned}
 CB(\text{PFC PV} \longrightarrow \text{PFC error x rule}) &= \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC error x rule1}) \\
 &+ \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC error x rule2}) \\
 &- \bar{W}(\text{PFC PV rule1} \longrightarrow \text{PFC error x rule2}) \\
 &- \bar{W}(\text{PFC PV rule2} \longrightarrow \text{PFC error x rule1})
 \end{aligned} \tag{24}$$

$$\begin{aligned}
 CB(\text{PFC error x rule} \longrightarrow \text{PFC E}) &= \bar{W}(\text{PFC error x rule1} \longrightarrow \text{PFC E rule1}) \\
 &+ \bar{W}(\text{PFC error x rule2} \longrightarrow \text{PFC E rule2}) \\
 &- \bar{W}(\text{PFC error x rule1} \longrightarrow \text{PFC E rule2}) \\
 &- \bar{W}(\text{PFC error x rule2} \longrightarrow \text{PFC E rule1})
 \end{aligned} \tag{25}$$

$$\begin{aligned}
 CB(\text{PFC error x rule} \longrightarrow \text{PFC PV}) &= \bar{W}(\text{PFC error x rule1} \longrightarrow \text{PFC PV rule1}) \\
 &+ \bar{W}(\text{PFC error x rule2} \longrightarrow \text{PFC PV rule2}) \\
 &- \bar{W}(\text{PFC error x rule1} \longrightarrow \text{PFC PV rule2}) \\
 &- \bar{W}(\text{PFC error x rule2} \longrightarrow \text{PFC PV rule1})
 \end{aligned} \tag{26}$$

The connectivity biases between the different response location-selective populations in the sensorimotor module (SM) are defined as

$$\begin{aligned}
 CB(\text{SM response E} \longrightarrow \text{SM response E}) &= \bar{W}(\text{SM E response 1} \longrightarrow \text{SM E response 1}) \\
 &+ \bar{W}(\text{SM E response 2} \longrightarrow \text{SM E response 2}) \\
 &+ \bar{W}(\text{SM E response 3} \longrightarrow \text{SM E response 3}) \\
 &- \bar{W}(\text{SM E response 1} \longrightarrow \text{SM E response 2 and 3}) \\
 &- \bar{W}(\text{SM E response 2} \longrightarrow \text{SM E response 1 and 3}) \\
 &- \bar{W}(\text{SM E response 3} \longrightarrow \text{SM E response 1 and 2}).
 \end{aligned} \tag{27}$$

In the last equation, for example, $\bar{W}(\text{SM response 1} \longrightarrow \text{SM response 2 and 3})$ represents the mean connection strength from excitatory neurons in the sensorimotor module that prefer response location 1 to those that prefer response locations 2 and 3.

766 The other connectivity biases were defined similarly

$$\begin{aligned}
 CB(\text{SM response E} \longrightarrow \text{SM response PV}) &= \bar{W}(\text{SM E response 1} \longrightarrow \text{SM PV response 1}) \\
 &+ \bar{W}(\text{SM E response 2} \longrightarrow \text{SM PV response 2}) \\
 &+ \bar{W}(\text{SM E response 3} \longrightarrow \text{SM PV response 3}) \\
 &- \bar{W}(\text{SM E response 1} \longrightarrow \text{SM PV response 2 and 3}) \\
 &- \bar{W}(\text{SM E response 2} \longrightarrow \text{SM PV response 1 and 3}) \\
 &- \bar{W}(\text{SM E response 3} \longrightarrow \text{SM PV response 1 and 2}).
 \end{aligned}
 \tag{28}$$

767

$$\begin{aligned}
 CB(\text{SM response PV} \longrightarrow \text{SM response E}) &= \bar{W}(\text{SM PV response 1} \longrightarrow \text{SM E response 2 and 3}) \\
 &+ \bar{W}(\text{SM PV response 2} \longrightarrow \text{SM E response 1 and 3}) \\
 &+ \bar{W}(\text{SM PV response 3} \longrightarrow \text{SM E response 1 and 2}) \\
 &- \bar{W}(\text{SM PV response 1} \longrightarrow \text{SM E response 1}) \\
 &- \bar{W}(\text{SM E response 2} \longrightarrow \text{SM PV response 2}) \\
 &- \bar{W}(\text{SM E response 3} \longrightarrow \text{SM PV response 3}).
 \end{aligned}
 \tag{29}$$

768

$$\begin{aligned}
 CB(\text{SM response PV} \longrightarrow \text{SM response PV}) &= \bar{W}(\text{SM PV response 1} \longrightarrow \text{SM PV response 2 and 3}) \\
 &+ \bar{W}(\text{SM PV response 2} \longrightarrow \text{SM PV response 1 and 3}) \\
 &+ \bar{W}(\text{SM PV response 3} \longrightarrow \text{SM PV response 1 and 2}) \\
 &- \bar{W}(\text{SM PV response 1} \longrightarrow \text{SM PV response 1}) \\
 &- \bar{W}(\text{SM PV response 2} \longrightarrow \text{SM PV response 2}) \\
 &- \bar{W}(\text{SM PV response 3} \longrightarrow \text{SM PV response 3}).
 \end{aligned}
 \tag{30}$$

769

770 The connectivity biases between the different rule-selective populations in the senso-
 771 rimotor module are defined as

$$\begin{aligned}
 CB(\text{SM rule E} \longrightarrow \text{SM rule E}) &= \bar{W}(\text{SM E rule 1} \longrightarrow \text{SM E rule 1}) \\
 &+ \bar{W}(\text{SM E rule 2} \longrightarrow \text{SM E rule 2}) \\
 &- \bar{W}(\text{SM E rule 1} \longrightarrow \text{SM E rule 2}) \\
 &- \bar{W}(\text{SM E rule 2} \longrightarrow \text{SM E rule 1}).
 \end{aligned}
 \tag{31}$$

772

773 The other connectivity biases were defined similarly

$$\begin{aligned}
 CB(\text{SM rule E} \rightarrow \text{SM rule PV}) &= \bar{W}(\text{SM E rule 1} \rightarrow \text{SM PV rule 1}) \\
 &+ \bar{W}(\text{SM E rule 2} \rightarrow \text{SM PV rule 2}) \\
 &- \bar{W}(\text{SM E rule 1} \rightarrow \text{SM PV rule 2}) \\
 &- \bar{W}(\text{SM E rule 2} \rightarrow \text{SM PV rule 1}).
 \end{aligned} \tag{32}$$

$$\begin{aligned}
 CB(\text{SM rule PV} \rightarrow \text{SM rule E}) &= \bar{W}(\text{SM PV rule 1} \rightarrow \text{SM E rule 2}) \\
 &+ \bar{W}(\text{SM PV rule 2} \rightarrow \text{SM E rule 1}) \\
 &- \bar{W}(\text{SM PV rule 1} \rightarrow \text{SM E rule 1}) \\
 &- \bar{W}(\text{SM PV rule 2} \rightarrow \text{SM E rule 2}).
 \end{aligned} \tag{33}$$

$$\begin{aligned}
 CB(\text{SM rule PV} \rightarrow \text{SM rule PV}) &= \bar{W}(\text{SM PV rule 1} \rightarrow \text{SM PV rule 2}) \\
 &+ \bar{W}(\text{SM PV rule 2} \rightarrow \text{SM PV rule 1}) \\
 &- \bar{W}(\text{SM PV rule 1} \rightarrow \text{SM PV rule 1}) \\
 &- \bar{W}(\text{SM PV rule 2} \rightarrow \text{SM PV rule 2}).
 \end{aligned} \tag{34}$$

777 The connectivity biases between the different shared feature-selective populations in
 778 the sensorimotor module are defined similarly. For the populations selective for the two
 779 colors

$$\begin{aligned}
 CB(\text{SM share feature (color) E} \rightarrow \text{SM shared feature (color) E}) &= \bar{W}(\text{SM E blue} \rightarrow \text{SM E blue}) \\
 &+ \bar{W}(\text{SM E red} \rightarrow \text{SM E red}) \\
 &- \bar{W}(\text{SM E blue} \rightarrow \text{SM E red}) \\
 &- \bar{W}(\text{SM E red} \rightarrow \text{SM E blue}),
 \end{aligned} \tag{35}$$

780 where for example $\bar{W}(\text{SM E blue} \rightarrow \text{SM E blue})$ is the average connection strength
 781 within the neural population selective for the shared feature blue.
 782

783 The other connectivity biases were defined similarly

$$\begin{aligned}
 CB(\text{SM shared feature (color) E} \rightarrow \text{SM shared feature (color) PV}) &= \bar{W}(\text{SM E blue} \rightarrow \text{SM PV blue}) \\
 &+ \bar{W}(\text{SM E red} \rightarrow \text{SM PV red}) \\
 &- \bar{W}(\text{SM E blue} \rightarrow \text{SM PV red}) \\
 &- \bar{W}(\text{SM E red} \rightarrow \text{SM PV blue}).
 \end{aligned} \tag{36}$$

784

$$\begin{aligned} CB(\text{SM shared feature (color) PV} \longrightarrow \text{SM shared feature (color) E}) &= \bar{W}(\text{SM PV blue} \longrightarrow \text{SM E red}) \\ &+ \bar{W}(\text{SM PV red} \longrightarrow \text{SM E blue}) \\ &- \bar{W}(\text{SM PV blue} \longrightarrow \text{SM E blue}) \\ &- \bar{W}(\text{SM PV red} \longrightarrow \text{SM E red}). \end{aligned} \tag{37}$$

785

$$\begin{aligned} CB(\text{SM shared feature (color) PV} \longrightarrow \text{SM shared feature (color) PV}) &= \bar{W}(\text{SM PV blue} \longrightarrow \text{SM PV red}) \\ &+ \bar{W}(\text{SM PV red} \longrightarrow \text{SM PV blue}) \\ &- \bar{W}(\text{SM PV blue} \longrightarrow \text{SM PV blue}) \\ &- \bar{W}(\text{SM PV red} \longrightarrow \text{SM PV red}). \end{aligned} \tag{38}$$

786

787 The connectivity biases between populations selective for different shared shapes were
788 defined analogously by substituting blue and red with circle and triangle.

789 **Simulation of the optogenetic inhibition**

790 Optogenetic inhibition was simulated by clamping the activity of neurons at 0
791 throughout the entire trial and the inter-trial interval.

792 **Principal angle between subspaces**

793 The principal angle between two subspaces is a generalization of angle between lines
794 and planes in Euclidean space to arbitrary dimensions [86]. It can be computed by iteratively
795 finding pairs of unit length “principal vectors”, one from each subspace, that have the
796 greatest inner product, subject to the condition that the principal vectors are orthogonal
797 to all previous principal vectors [87].

798 In computing the principal angles between different rule-selective and response-
799 selective subspaces, we first determined the dimensionality of the subspaces using the par-
800 ticipation ratio [88]. Then the principal angles were computed using the “subspace_angles”
801 function from the Python package Scipy. The largest principal angle was used.

802 To obtain a shuffled distribution, we first evenly split all trials belonging to a particular
803 rule or response into two halves. Then, we generated two subspaces from neural trajectories
804 during the two group of trials. A principal angle between these two subspaces was then
805 computed for each rule/response. The angles were then averaged across all rules/responses

806 to obtain a principal angle from shuffled data. This process was repeated 100 times to
807 generate a distribution of principal angles from shuffled data.

808 **Assessing the strength of non-linear mixed selectivity**

The extent to which neurons in the sensorimotor module encode the conjunction of stimulus and rule in a non-linear fashion was evaluated using the coefficient of determination of a linear regression model. To tease apart non-linear and linear mixed selectivity, we first fitted the mean activity of each neuron during response period using a set of regressors that represent either the rule or the stimulus alone:

$$FR(n, tr) = \sum_s \beta_{n,s} \mathbb{1}(\text{stim}(tr) = s) + \sum_r \beta_{n,r} \mathbb{1}(\text{rule}(tr) = r), \quad (39)$$

809 where $FR(n, tr)$ is the firing rate of neuron n during trial tr . $\mathbb{1}$ is the indicator function.
810 For example, $\mathbb{1}(\text{stim}(tr) = s) = 1$ if the stimulus during trial tr is s , and it is 0 otherwise.

Then, another linear regression model was fitted on the residual activity unexplained by the linear regression model above, using the conjunction of rule and stimulus as regressors:

$$\tilde{FR}(n, tr) = \sum_{s,r} \beta_{n,s,r} \mathbb{1}(\text{stim}(tr) = s, \text{rule}(tr) = r), \quad (40)$$

811 where $\tilde{FR}(n, tr)$ is the firing rate of neuron n during trial tr subtracted by the pre-
812 dicted firing rate from the model defined by Equation 39. The R^2 value of this regression
813 model was used to represent the strength of non-linear mixed selectivity.

References

814

- 815 [1] David A Grant and Esta Berg. A behavioral analysis of degree of reinforcement and
816 ease of shifting to new responses in a weigl-type card-sorting problem. *Journal of*
817 *experimental psychology*, 38(4):404, 1948.
- 818 [2] Brenda Milner. Effects of different brain lesions on card sorting: The role of the frontal
819 lobes. *Archives of neurology*, 9(1):90–100, 1963.
- 820 [3] R Dias, TW Robbins, and AC Roberts. Primate analogue of the wisconsin card sorting
821 test: effects of excitotoxic lesions of the prefrontal cortex in the marmoset. *Behavioral*
822 *neuroscience*, 110:872, 1996.
- 823 [4] RE Passingham. Non-reversal shifts after selective prefrontal ablations in monkeys
824 (macaca mulatta). *Neuropsychologia*, 10(1):41–46, 1972.
- 825 [5] Katsuyuki Sakai. Task set and prefrontal cortex. *Annu. Rev. Neurosci.*, 31:219–245,
826 2008.
- 827 [6] M. J. Buckley, F. A. Mansouri, H. Hoda, M. Mahboubi, P. G. Browning, S. C. Kwok,
828 A. Phillips, and K. Tanaka. Dissociable components of rule-guided behavior depend
829 on distinct medial and prefrontal regions. *Science*, 325:52–58, 2009.
- 830 [7] E. K. Miller and J. D. Cohen. An integrative theory of prefrontal cortex function.
831 *Annual Review of Neuroscience*, 24:167–202, 2001.
- 832 [8] Farshad A Mansouri, Kenji Matsumoto, and Keiji Tanaka. Prefrontal cell activities
833 related to monkeys’ success and failure in adapting to rule changes in a wisconsin card
834 sorting test analog. *Journal of Neuroscience*, 26(10):2745–2756, 2006.
- 835 [9] Tsukasa Kamigaki, Tetsuya Fukushima, and Yasushi Miyashita. Cognitive set recon-
836 figuration signaled by macaque posterior parietal neurons. *Neuron*, 61(6):941–951,
837 2009.
- 838 [10] Morteza Sarafyazd and Mehrdad Jazayeri. Hierarchical reasoning by neural circuits in
839 the frontal cortex. *Science*, 364(6441):eaav8911, 2019.
- 840 [11] Takuya Ito, Guangyu Robert Yang, Patryk Laurent, Douglas H Schultz, and Michael W
841 Cole. Constructing neural network models from brain data reveals representational
842 transformations linked to adaptive behavior. *Nature communications*, 13(1):673, 2022.
- 843 [12] Kristen Delevich, Jason Tucciarone, Z Josh Huang, and Bo Li. The mediodorsal tha-
844 lamus drives feedforward inhibition in the anterior cingulate cortex via parvalbumin
845 interneurons. *Journal of Neuroscience*, 35(14):5743–5753, 2015.

- 846 [13] Hyun-Jae Pi, Balázs Hangya, Duda Kvitsiani, Joshua I Sanders, Z Josh Huang, and
847 Adam Kepecs. Cortical interneurons that specialize in disinhibitory control. *Nature*,
848 503(7477):521–524, 2013.
- 849 [14] Siyu Zhang, Min Xu, Tsukasa Kamigaki, Johnny Phong Hoang Do, Wei-Cheng Chang,
850 Sean Jenvay, Kazunari Miyamichi, Liqun Luo, and Yang Dan. Long-range and local
851 circuits for top-down modulation of visual cortex processing. *Science*, 345(6197):660–
852 665, 2014.
- 853 [15] William Muñoz, Robin Tremblay, Daniel Levenstein, and Bernardo Rudy. Layer-
854 specific modulation of neocortical dendritic inhibition during active wakefulness. *Sci-
855 ence*, 355(6328):954–959, 2017.
- 856 [16] Andreas J Keller, Mario Dipoppa, Morgane M Roth, Matthew S Caudill, Alessan-
857 dro Ingrassia, Kenneth D Miller, and Massimo Scanziani. A disinhibitory circuit for
858 contextual modulation in primary visual cortex. *Neuron*, 108(6):1181–1193, 2020.
- 859 [17] X. J. Wang, J. Tegnér, C. Constantinidis, and P. S. Goldman-Rakic. Division of labor
860 among distinct subtypes of inhibitory neurons in a cortical microcircuit of working
861 memory. *Proceedings of the National Academy of Science, USA*, 101(5):1368–73, 2004.
- 862 [18] A. Kepecs and G. Fishell. Interneuron cell types are fit to function. *Nature*, 505:318–
863 326, 2014.
- 864 [19] R. Tremblay, S. Lee, and B. Rudy. GABAergic interneurons in the neocortex: From
865 cellular properties to circuits. *Neuron*, 91:260–292, 2016.
- 866 [20] Guangyu Robert Yang, John D Murray, and Xiao-Jing Wang. A dendritic disinhibitory
867 circuit mechanism for pathway-specific gating. *Nature communications*, 7(1):1–14,
868 2016.
- 869 [21] Kiyoshi Nakahara, Toshihiro Hayashi, Seiki Konishi, and Yasushi Miyashita. Func-
870 tional mri of macaque monkeys performing a cognitive set-shifting task. *Science*,
871 295(5559):1532–1536, 2002.
- 872 [22] G. R. Yang and X.-J. Wang. Artificial neural networks for neuroscientists: a primer.
873 *Neuron*, 107:1048–1070, 2020.
- 874 [23] Monika Jadi, Alon Polsky, Jackie Schiller, and Bartlett W Mel. Location-dependent
875 effects of inhibition on local spiking in pyramidal neuron dendrites. *PLoS computational
876 biology*, 8(6):e1002550, 2012.

- 877 [24] Valerio Mante, David Sussillo, Krishna V Shenoy, and William T Newsome.
878 Context-dependent computation by recurrent dynamics in prefrontal cortex. *nature*,
879 503(7474):78–84, 2013.
- 880 [25] Guangyu Robert Yang, Madhura R Joglekar, H Francis Song, William T Newsome,
881 and Xiao-Jing Wang. Task representations in neural networks trained to perform many
882 cognitive tasks. *Nature neuroscience*, 22(2):297–306, 2019.
- 883 [26] H Francis Song, Guangyu R Yang, and Xiao-Jing Wang. Training excitatory-inhibitory
884 recurrent neural networks for cognitive tasks: a simple and flexible framework. *PLoS*
885 *computational biology*, 12(2):e1004792, 2016.
- 886 [27] Tsukasa Kamigaki, Tetsuya Fukushima, and Yasushi Miyashita. Neuronal signal dy-
887 namics during preparation and execution for behavioral shifting in macaque posterior
888 parietal cortex. *Journal of cognitive neuroscience*, 23(9):2503–2520, 2011.
- 889 [28] Joaquin M Fuster. Unit activity in prefrontal cortex during delayed-response perfor-
890 mance: neuronal correlates of transient memory. *Journal of neurophysiology*, 36(1):61–
891 78, 1973.
- 892 [29] Shintaro Funahashi, Charles J Bruce, and Patricia S Goldman-Rakic. Mnemonic coding
893 of visual space in the monkey’s dorsolateral prefrontal cortex. *Journal of neurophysi-*
894 *ology*, 61(2):331–349, 1989.
- 895 [30] P. Goldman-Rakic. Cellular basis of working memory. *Neuron*, 14:477–85, 1995.
- 896 [31] Ranulfo Romo, Carlos D Brody, Adrián Hernández, and Luis Lemus. Neuronal corre-
897 lates of parametric working memory in the prefrontal cortex. *Nature*, 399(6735):470–
898 473, 1999.
- 899 [32] Thomas B Christophel, P Christiaan Klink, Bernhard Spitzer, Pieter R Roelfsema, and
900 John-Dylan Haynes. The distributed nature of working memory. *Trends in cognitive*
901 *sciences*, 21(2):111–124, 2017.
- 902 [33] Matthew L Leavitt, Diego Mendoza-Halliday, and Julio C Martinez-Trujillo. Sustained
903 activity encoding working memories: not fully distributed. *Trends in Neurosciences*,
904 40(6):328–346, 2017.
- 905 [34] Zengcai V Guo, Hidehiko K Inagaki, Kayvon Daie, Shaul Druckmann, Charles R Ger-
906 fen, and Karel Svoboda. Maintenance of persistent activity in a frontal thalamocortical
907 loop. *Nature*, 545(7653):181–186, 2017.

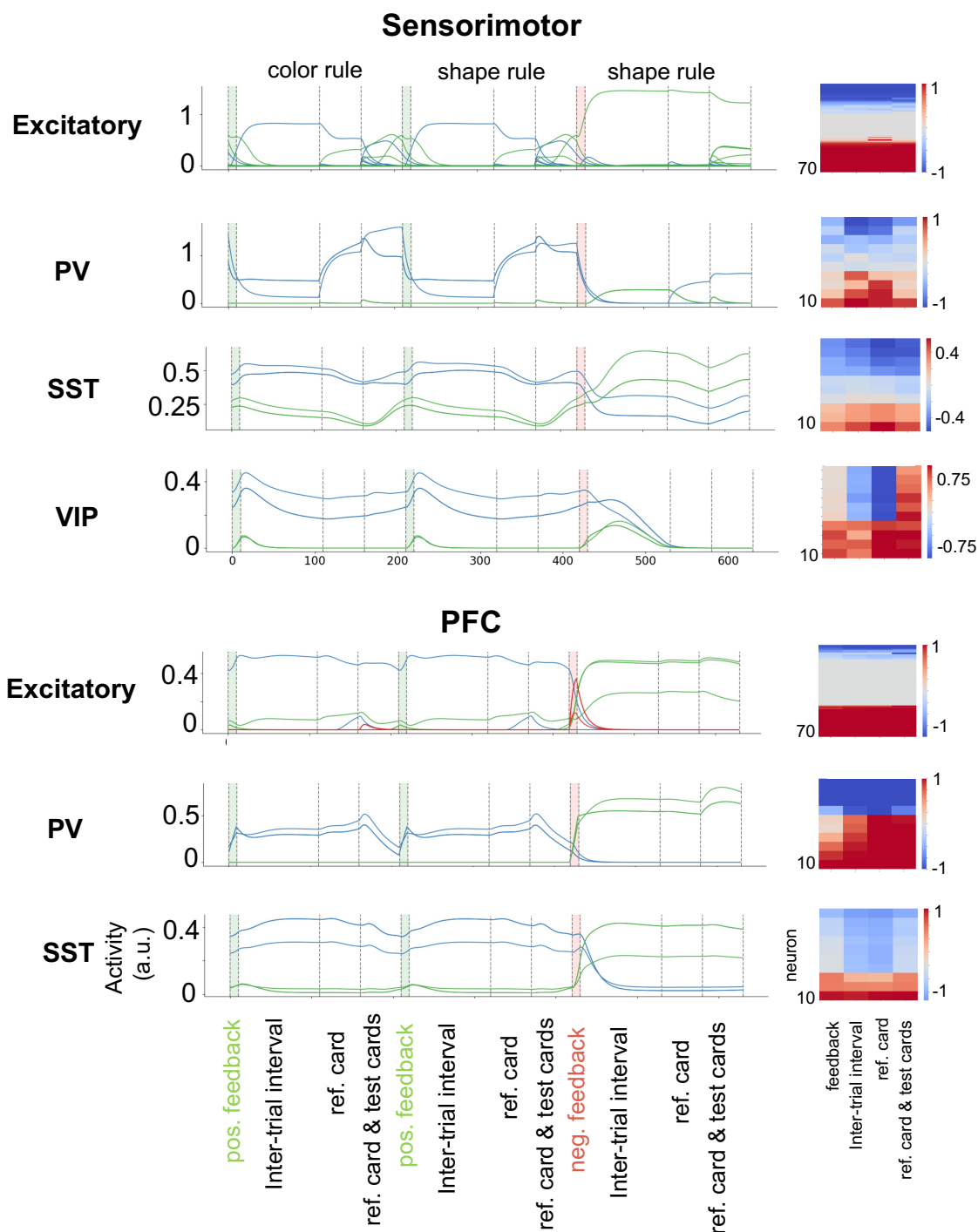
- 908 [35] Kartik K Sreenivasan and Mark D’Esposito. The what, where and how of delay activity.
909 *Nature reviews neuroscience*, 20(8):466–481, 2019.
- 910 [36] Kong-Fatt Wong and Xiao-Jing Wang. A recurrent network mechanism of time inte-
911 gration in perceptual decisions. *Journal of Neuroscience*, 26(4):1314–1328, 2006.
- 912 [37] Farzaneh Najafi, Gamaleldin F Elsayed, Robin Cao, Eftychios Pnevmatikakis, Peter E
913 Latham, John P Cunningham, and Anne K Churchland. Excitatory and inhibitory
914 subnetworks are equally selective during decision-making and emerge simultaneously
915 during learning. *Neuron*, 105(1):165–179, 2020.
- 916 [38] James P Roach, Anne K Churchland, and Tatiana A Engel. Choice selective inhi-
917 bition drives stability and competition in decision circuits. *Nature Communications*,
918 14(1):147, 2023.
- 919 [39] Hongbo Jia, Nathalie L Rochefort, Xiaowei Chen, and Arthur Konnerth. Dendritic
920 organization of sensory input to cortical neurons in vivo. *Nature*, 464(7293):1307–1312,
921 2010.
- 922 [40] Joseph Cichon and Wen-Biao Gan. Branch-specific dendritic ca²⁺ spikes cause persis-
923 tent synaptic plasticity. *Nature*, 520(7546):180–185, 2015.
- 924 [41] Shannon K Rashid, Victor Pedrosa, Martial A Dufour, Jason J Moore, Spyridon
925 Chavlis, Rodrigo G Delatorre, Panayiota Poirazi, Claudia Clopath, and Jayeeta Basu.
926 The dendritic spatial code: branch-specific place tuning and its experience-dependent
927 decoupling. *BioRxiv*, pages 2020–01, 2020.
- 928 [42] Jakob Voigts and Mark T Harnett. Somatic and dendritic encoding of spatial variables
929 in retrosplenial cortex differs during 2d navigation. *Neuron*, 105(2):237–245, 2020.
- 930 [43] Mattia Rigotti, Omri Barak, Melissa R Warden, Xiao-Jing Wang, Nathaniel D Daw,
931 Earl K Miller, and Stefano Fusi. The importance of mixed selectivity in complex
932 cognitive tasks. *Nature*, 497(7451):585–90, 2013.
- 933 [44] Atsushi Kikumoto and Ulrich Mayr. Conjunctive representations that integrate stim-
934 uli, responses, and rules are critical for action selection. *Proceedings of the National
935 Academy of Sciences*, 117(19):10603–10608, 2020.
- 936 [45] Atsushi Kikumoto, Ulrich Mayr, and David Badre. The role of conjunctive represen-
937 tations in prioritizing and selecting planned actions. *Elife*, 11:e80153, 2022.
- 938 [46] Jonathan D Wallis, Kathleen C Anderson, and Earl K Miller. Single neurons in pre-
939 frontal cortex encode abstract rules. *Nature*, 411(6840):953–956, 2001.

- 940 [47] Matthew M Botvinick, Jonathan D Cohen, and Cameron S Carter. Conflict monitoring
941 and anterior cingulate cortex: an update. *Trends in cognitive sciences*, 8(12):539–546,
942 2004.
- 943 [48] René Quilodran, Marie Rothe, and Emmanuel Procyk. Behavioral shifts and action
944 valuation in the anterior cingulate cortex. *Neuron*, 57(2):314–325, 2008.
- 945 [49] Nils Kolling, Marco K Wittmann, Tim EJ Behrens, Eerie D Boorman, Rogier B Mars,
946 and Matthew FS Rushworth. Value, search, persistence and model updating in anterior
947 cingulate cortex. *Nature neuroscience*, 19(10):1280–1285, 2016.
- 948 [50] Farshad Alizadeh Mansouri, David J Freedman, and Mark J Buckley. Emergence of
949 abstract rules in the primate brain. *Nature Reviews Neuroscience*, 21(11):595–610,
950 2020.
- 951 [51] Timothy Spellman, Malka Svei, Jesse Kaminsky, Gabriela Manzano-Nieves, and Conor
952 Liston. Prefrontal deep projection neurons enable cognitive flexibility via persistent
953 feedback monitoring. *Cell*, 184(10):2750–2766, 2021.
- 954 [52] C Daniel Salzman and Stefano Fusi. Emotion, cognition, and mental state represen-
955 tation in amygdala and prefrontal cortex. *Annual review of neuroscience*, 33:173–202,
956 2010.
- 957 [53] Clay B Holroyd and Michael GH Coles. The neural basis of human error processing: re-
958 inforcement learning, dopamine, and the error-related negativity. *Psychological review*,
959 109(4):679, 2002.
- 960 [54] Masayuki Matsumoto and Okihide Hikosaka. Two types of dopamine neuron distinctly
961 convey positive and negative motivational signals. *Nature*, 459(7248):837–841, 2009.
- 962 [55] Richard A Andersen and He Cui. Intention, action planning, and decision making in
963 parietal-frontal circuits. *Neuron*, 63(5):568–583, 2009.
- 964 [56] Bernard W Balleine and John P O’doherly. Human and rodent homologues in action
965 control: corticostriatal determinants of goal-directed and habitual action. *Neuropsy-
966 chopharmacology*, 35(1):48–69, 2010.
- 967 [57] Bevil R Conway. Color vision, cones, and color-coding in the cortex. *The neuroscientist*,
968 15(3):274–290, 2009.
- 969 [58] Rosa Lafer-Sousa and Bevil R Conway. Parallel, multi-stage processing of colors, faces
970 and shapes in macaque inferior temporal cortex. *Nature neuroscience*, 16(12):1870–
971 1878, 2013.

- 972 [59] Le Chang, Pinglei Bao, and Doris Y Tsao. The representation of colored objects in
973 macaque color patches. *Nature communications*, 8(1):2064, 2017.
- 974 [60] Mark M Churchland, John P Cunningham, Matthew T Kaufman, Justin D Foster,
975 Paul Nuyujukian, Stephen I Ryu, and Krishna V Shenoy. Neural population dynamics
976 during reaching. *Nature*, 487(7405):51–56, 2012.
- 977 [61] Nicholas A Steinmetz, Peter Zatzka-Haas, Matteo Carandini, and Kenneth D Harris.
978 Distributed coding of choice, action and engagement across the mouse brain. *Nature*,
979 576(7786):266–273, 2019.
- 980 [62] Shih-Yi Tseng, Selmaan N Chettih, Charlotte Arlt, Roberto Barroso-Luque, and
981 Christopher D Harvey. Shared and specialized coding across posterior cortical areas
982 for dynamic navigation decisions. *Neuron*, 110(15):2484–2502, 2022.
- 983 [63] Charles Findling, Felix Hubert, International Brain Laboratory, Luigi Acerbi, Brandon
984 Benson, Julius Benson, Daniel Birman, Niccolò Bonacchi, Matteo Carandini, Joana A
985 Catarino, et al. Brain-wide representations of prior information in mouse decision-
986 making. *BioRxiv*, pages 2023–07, 2023.
- 987 [64] International Brain Lab, Brandon Benson, Julius Benson, Daniel Birman, Niccolò
988 Bonacchi, Matteo Carandini, Joana A Catarino, Gaelle A Chapuis, Anne K Church-
989 land, Yang Dan, et al. A brain-wide map of neural activity during complex behaviour.
990 *bioRxiv*, pages 2023–07, 2023.
- 991 [65] Sean Froudish-Walsh, Daniel P Bliss, Xingyu Ding, Lucija Rapan, Meiqi Niu, Kenneth
992 Knoblauch, Karl Zilles, Henry Kennedy, Nicola Palomero-Gallagher, and Xiao-Jing
993 Wang. A dopamine gradient controls access to distributed working memory in the
994 large-scale monkey cortex. *Neuron*, 109(21):3500–3520, 2021.
- 995 [66] Jorge F Mejías and Xiao-Jing Wang. Mechanisms of distributed working memory in a
996 large-scale network of macaque neocortex. *Elife*, 11:e72136, 2022.
- 997 [67] Mattia Rigotti, Daniel Ben Dayan Rubin, Xiao-Jing Wang, and Stefano Fusi. Internal
998 representation of task rules by recurrent dynamics: the importance of the diversity of
999 neural responses. *Frontiers in computational neuroscience*, 4:24, 2010.
- 1000 [68] Stanislas Dehaene and Jean-Pierre Changeux. The wisconsin card sorting test: The-
1001 oretical analysis and modeling in a neuronal network. *Cerebral cortex*, 1(1):62–79,
1002 1991.

- 1003 [69] Daniel Turner-Evans, Stephanie Wegener, Herve Rouault, Romain Franconville, Tanya
1004 Wolff, Johannes D Seelig, Shaul Druckmann, and Vivek Jayaraman. Angular velocity
1005 integration in a fly heading circuit. *Elife*, 6:e23496, 2017.
- 1006 [70] Benjamin Y Hayden, John M Pearson, and Michael L Platt. Neuronal basis of sequen-
1007 tial foraging decisions in a patchy environment. *Nature neuroscience*, 14(7):933–939,
1008 2011.
- 1009 [71] João D Semedo, Amin Zandvakili, Christian K Machens, M Yu Byron, and Adam Kohn.
1010 Cortical areas interact through a communication subspace. *Neuron*, 102(1):249–259,
1011 2019.
- 1012 [72] Christopher Langdon, Mikhail Genkin, and Tatiana A Engel. A unifying perspective
1013 on neural manifolds and circuits for cognition. *Nature Reviews Neuroscience*, pages
1014 1–15, 2023.
- 1015 [73] Braden A Purcell and Roozbeh Kiani. Hierarchical decision processes that operate over
1016 distinct timescales underlie choice and changes in strategy. *Proceedings of the national
1017 academy of sciences*, 113(31):E4531–E4540, 2016.
- 1018 [74] Kevin Johnston, Helen M Levin, Michael J Koval, and Stefan Everling. Top-down
1019 control-signal dynamics in anterior cingulate and prefrontal cortex neurons following
1020 task switching. *Neuron*, 53(3):453–462, 2007.
- 1021 [75] Ken-Ichiro Tsutsui, Takayuki Hosokawa, Munekazu Yamada, and Toshio Iijima. Repre-
1022 sentation of functional category in the monkey prefrontal cortex and its rule-dependent
1023 use for behavioral selection. *Journal of Neuroscience*, 36(10):3038–3048, 2016.
- 1024 [76] Silvia Bernardi, Marcus K Benna, Mattia Rigotti, Jérôme Munuera, Stefano Fusi, and
1025 C Daniel Salzman. The geometry of abstraction in the hippocampus and prefrontal
1026 cortex. *Cell*, 183(4):954–967, 2020.
- 1027 [77] Vishwa Goudar, Jeong-Woo Kim, Yue Liu, Adam JO Dede, Michael J Jutras, Ivan
1028 Skelin, Michael Ruvalcaba, William Chang, Adrienne L Fairhall, Jack J Lin, et al.
1029 Comparing rapid rule-learning strategies in humans and monkeys. *bioRxiv*, pages 2023–
1030 01, 2023.
- 1031 [78] Steven W Kennerley, Mark E Walton, Timothy EJ Behrens, Mark J Buckley, and
1032 Matthew FS Rushworth. Optimal decision making and the anterior cingulate cortex.
1033 *Nature neuroscience*, 9(7):940–947, 2006.
- 1034 [79] Cheng Xue, Lily E Kramer, and Marlene R Cohen. Dynamic task-belief is an integral
1035 part of decision-making. *Neuron*, 110(15):2503–2511, 2022.

- 1036 [80] Ido Ben-Artzi, Yoav Kessler, Bruno Nicenboim, and Nitzan Shahar. Computational
1037 mechanisms underlying latent value updating of unchosen actions. *Science Advances*,
1038 9(42):eadi2704, 2023.
- 1039 [81] Amitai Shenhav, Matthew M Botvinick, and Jonathan D Cohen. The expected value of
1040 control: an integrative theory of anterior cingulate cortex function. *Neuron*, 79(2):217–
1041 240, 2013.
- 1042 [82] Stefano Fusi, Wael F Asaad, Earl K Miller, and Xiao-Jing Wang. A neural circuit
1043 model of flexible sensorimotor mapping: learning and forgetting on multiple timescales.
1044 *Neuron*, 54(2):319–333, 2007.
- 1045 [83] Xiaolong Jiang, Shan Shen, Cathryn R Cadwell, Philipp Berens, Fabian Sinz, Alexan-
1046 der S Ecker, Saumil Patel, and Andreas S Tolias. Principles of connectivity among
1047 morphologically defined cell types in adult neocortex. *Science*, 350(6264):aac9462,
1048 2015.
- 1049 [84] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Pro-
1050 ceedings of the IEEE*, 78(10):1550–1560, 1990.
- 1051 [85] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization.
1052 *arXiv preprint arXiv:1412.6980*, 2014.
- 1053 [86] Camille Jordan. Essai sur la géométrie à n dimensions. *Bulletin de la Société mathé-
1054 matique de France*, 3:103–174, 1875.
- 1055 [87] ke Björck and Gene H Golub. Numerical methods for computing angles between linear
1056 subspaces. *Mathematics of computation*, 27(123):579–594, 1973.
- 1057 [88] Peiran Gao, Eric Trautmann, Byron Yu, Gopal Santhanam, Stephen Ryu, Krishna
1058 Shenoy, and Surya Ganguli. A theory of multineuronal dimensionality, dynamics and
1059 measurement. *BioRxiv*, page 214262, 2017.

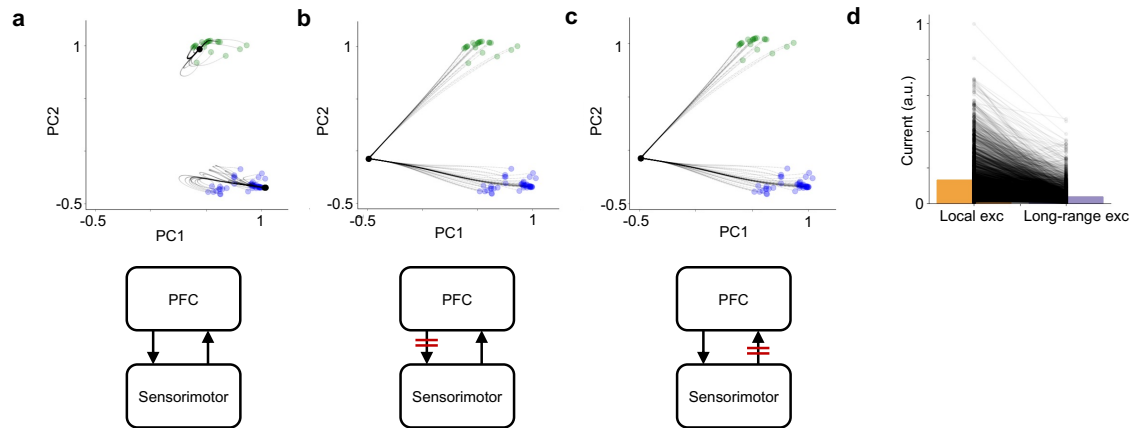


Supplementary Figure 1: Activity of single neurons.

Left: activity from example neurons across three consecutive trials with a rule switch. Blue and green traces represent neurons with higher activity during the color rule and shape rule respectively. Red traces for the PFC excitatory neurons represent neurons that are preferentially activated by negative feedback.

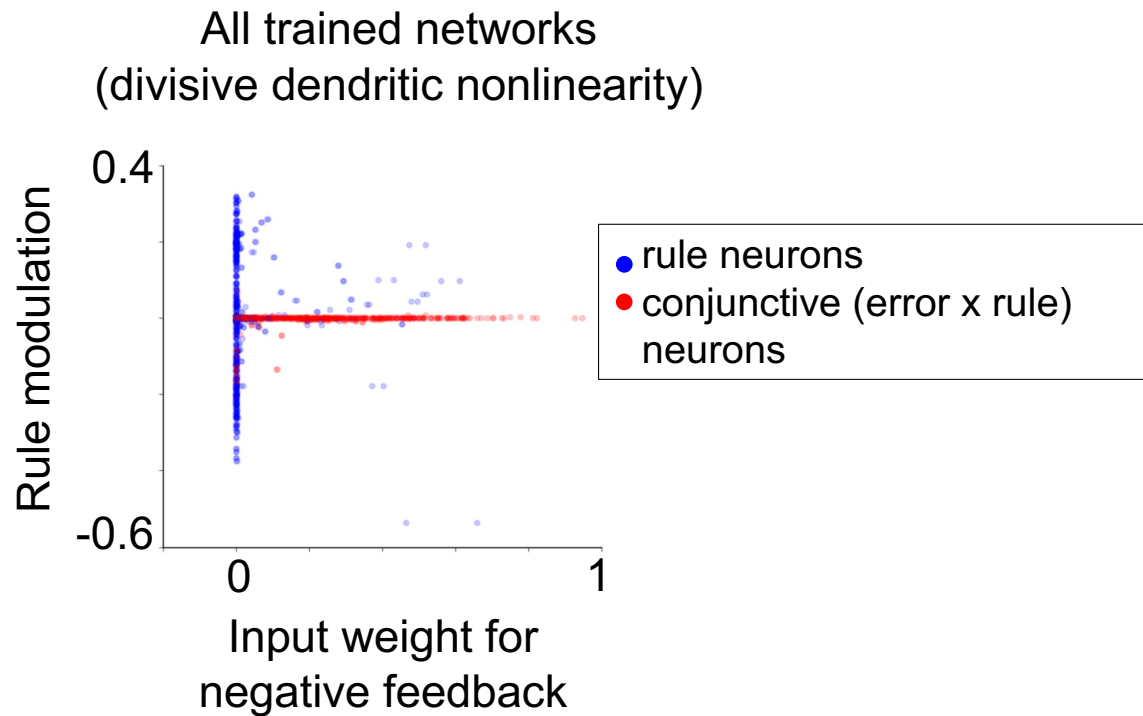
Right: Rule selectivity across neurons and task epochs.

VIP neurons in the PFC module do not receive any excitatory inputs, therefore are not shown.

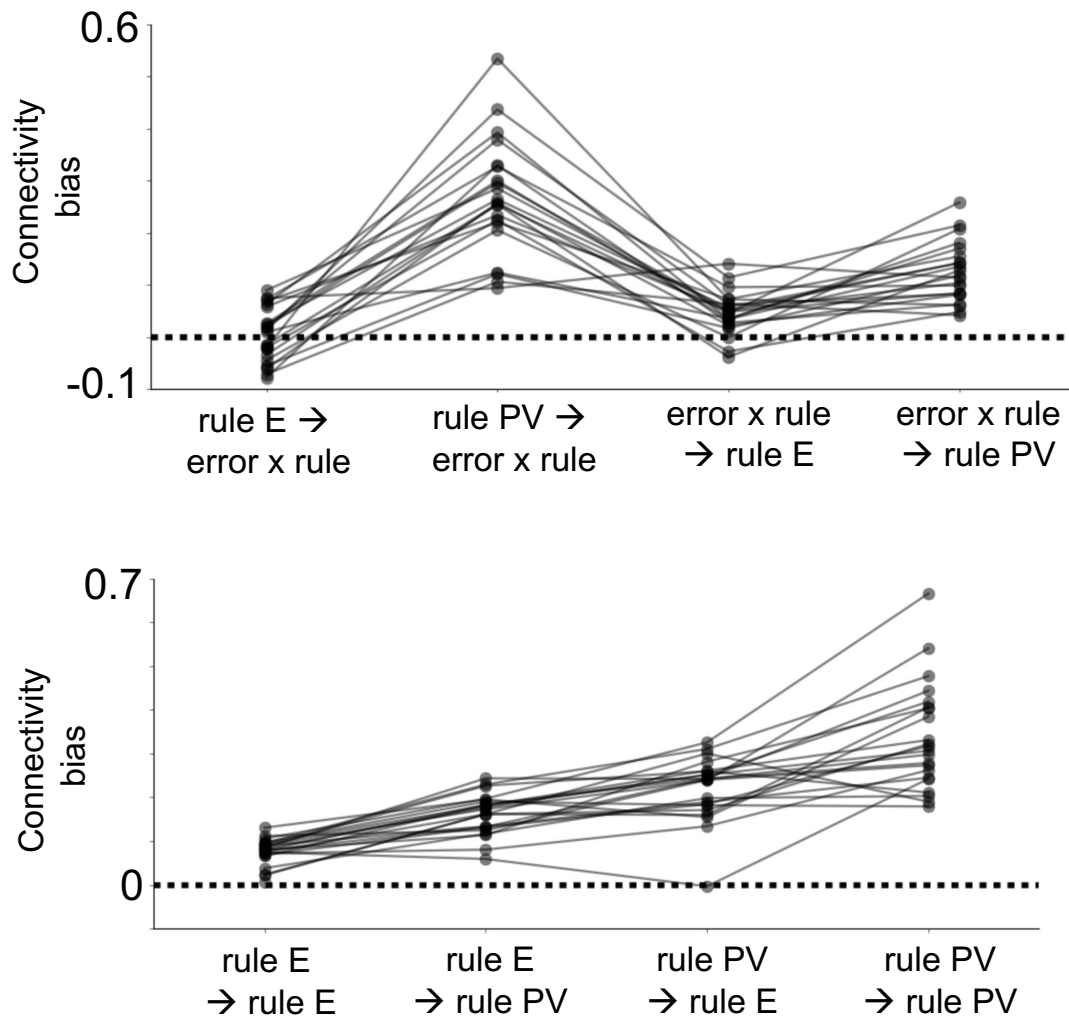


Supplementary Figure 2: Two attractor states in the PFC module supported by inter-modular connections.

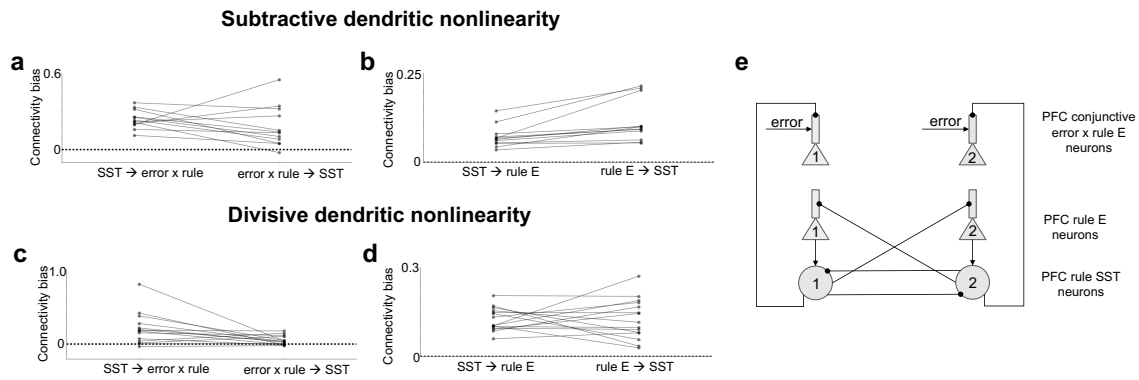
The neural trajectories that represent the autonomous dynamics of the trained RNN with intact inter-modular connections (a), when the connections from the PFC module to the sensorimotor module were lesioned (b), and when the connections from the sensorimotor module to the PFC module were lesioned (c). The networks started from random time point during color rule (blue) and shape rule (green) trials. Black points represent the final states of the networks after 500 timesteps (5 seconds). d. Excitatory and PV neurons in the PFC module receive stronger local excitation than long-range excitation. Each line represents a neuron. Bars represent average across neurons. Result aggregated across all trained networks. Student's t-test, $p < .001$.



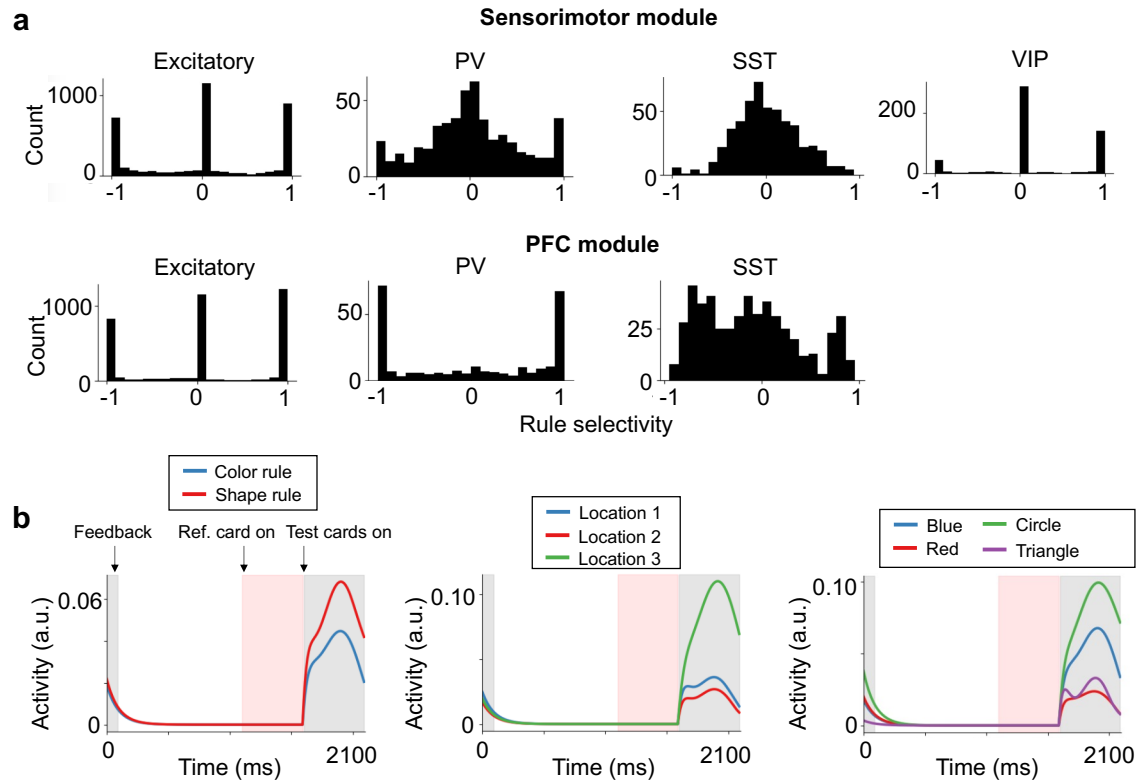
Supplementary Figure 3: The emergence of two populations of excitatory neurons in the PFC module of networks with divisive dendritic nonlinearity. The rule modulation against input weight for negative feedback for all the rule neurons and conjunctive error x rule neurons in the PFC module of networks with divisive dendritic nonlinearity (c.f. Figure 2c).



Supplementary Figure 4: The connectivity biases between different subpopulations of excitatory and PV neurons in the PFC module of networks with divisive dendritic nonlinearity (c.f. Figure 3d).



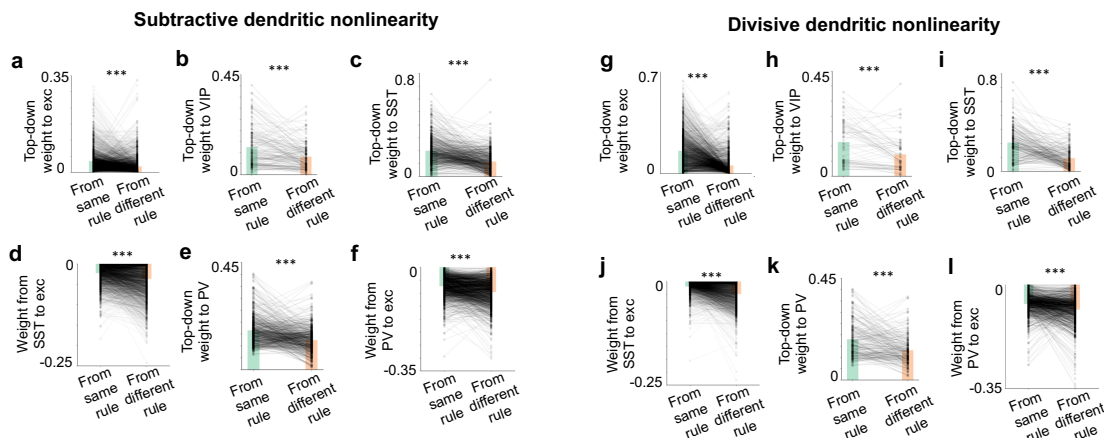
Supplementary Figure 5: The connectivity structure between the excitatory neurons and SST neurons in the PFC module. **a.** The connectivity bias between the SST neurons and the conjunctive error x rule neurons in the PFC module, for networks with subtractive dendritic nonlinearity. **b.** The connectivity bias between the SST neurons and the rule neurons in the PFC module, for networks with subtractive dendritic nonlinearity. **c - d.** The same connectivity biases but for networks with divisive dendritic nonlinearity. **e** The connectivity structure between the SST and excitatory neurons in the PFC module as revealed by **a- d**. Only the stronger connections are plotted. A similar connectivity pattern exists between SST and excitatory neurons as compared to the connectivity pattern between PV and excitatory neurons (c.f. Figure 3c).



Supplementary Figure 6: Single neurons feature selectivity in the PFC and sensorimotor modules

a. The distribution of rule selectivity across different cell types. Result is aggregated across all trained networks. Result for VIP neurons in the PFC module is not shown since they do not receive excitation.

b. The trial-averaged activity of an example excitatory neuron in the sensorimotor module with preferred rule $R = color\ rule$, response location $L = 3$ and shared feature $F = circle$. Trials were sorted according to rule (left), response (middle) and shared feature (right).



Supplementary Figure 7: Structure in the top-down projections, across all networks.

a. Each line represents the mean connection strength onto one excitatory neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent the mean across neurons. PFC excitatory neurons project more strongly to the excitatory neurons in the sensorimotor module that prefer the same rule (c.f. Figure 4a).

b. Each line represents the mean connection strength onto one VIP neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent the mean across neurons. PFC excitatory neurons project more strongly to the VIP neurons in the sensorimotor module that prefer the same rule (c.f. Figure 4b).

c. Each line represents the mean connection strength onto one PV neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent the mean across neurons. PFC excitatory neurons project more strongly to the PV neurons in the sensorimotor module that prefer the same rule (c.f. Figure 4c).

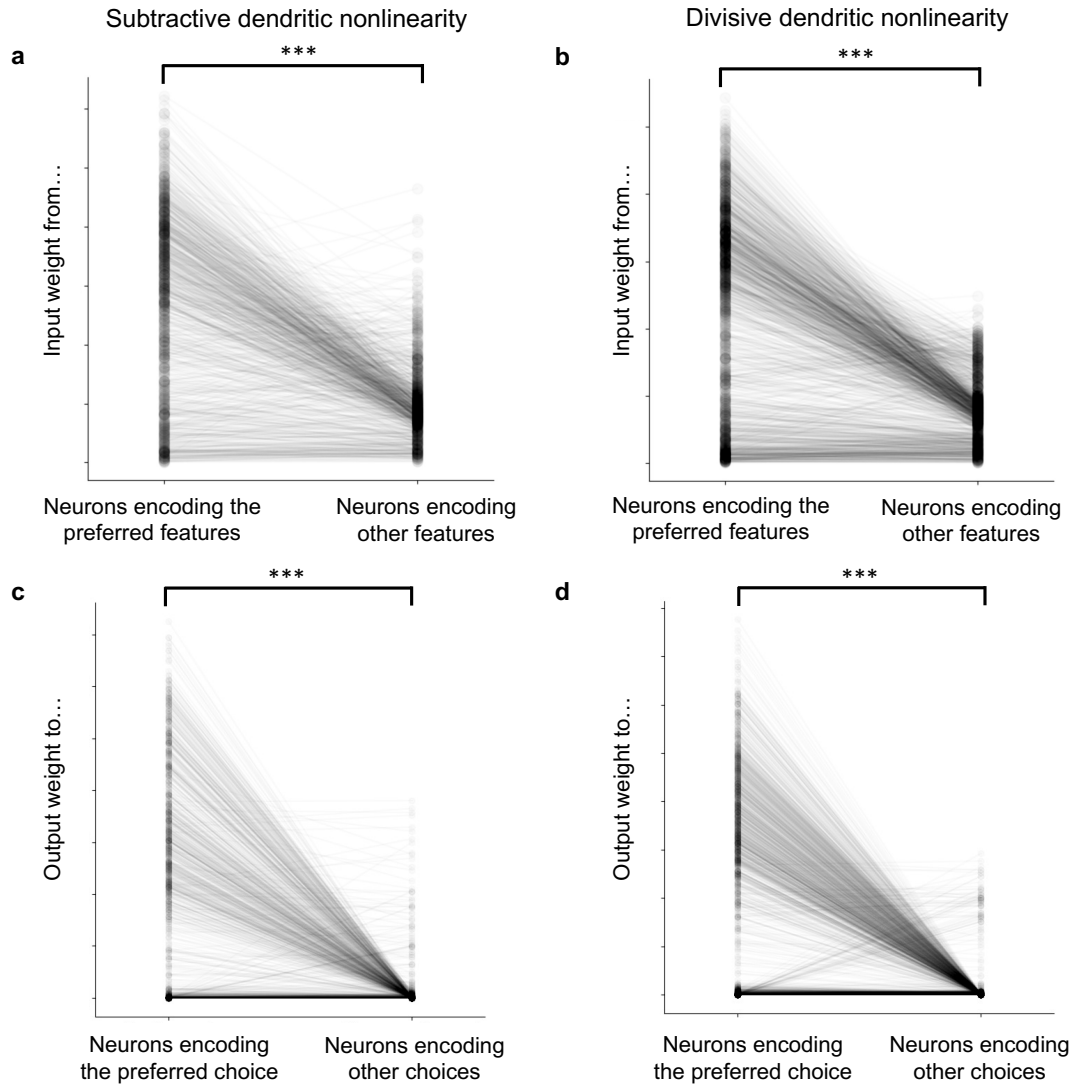
d. Each line represents the mean connection strength onto one SST neuron in the sensorimotor module, from the PFC excitatory neurons that prefer the same rule and the different rule. Bars represent the mean across neurons. PFC excitatory neurons project more strongly to local SST neurons that prefer the same rule (c.f. Figure 4d).

e. Each line represents the mean connection strength onto one excitatory neuron of the sensorimotor module, from the SST neurons in the sensorimotor module that prefer the same rule and the different rule. Bars represent mean across neurons. SST neurons in the sensorimotor module project more strongly to local excitatory neurons that prefer the opposite rule (c.f. Figure 4e).

f. Each line represents the mean connection strength onto one excitatory neuron of the sensorimotor module, from the PV neurons in the sensorimotor module that prefer the same rule and the different rule. Bars represent mean across neurons. PV neurons in the sensorimotor module project more strongly to local excitatory neurons that prefer the opposite rule (c.f. Figure 4f).

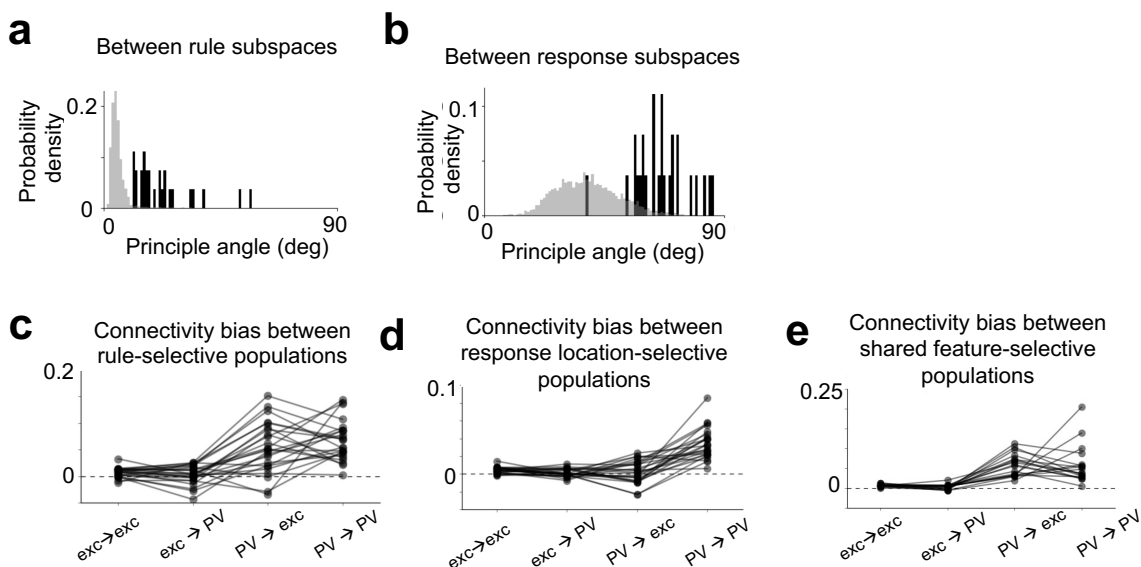
g - l. Same as **a - f** for networks with the divisive dendritic nonlinearity.

$p < 0.001$ for all panels, Student's t test.

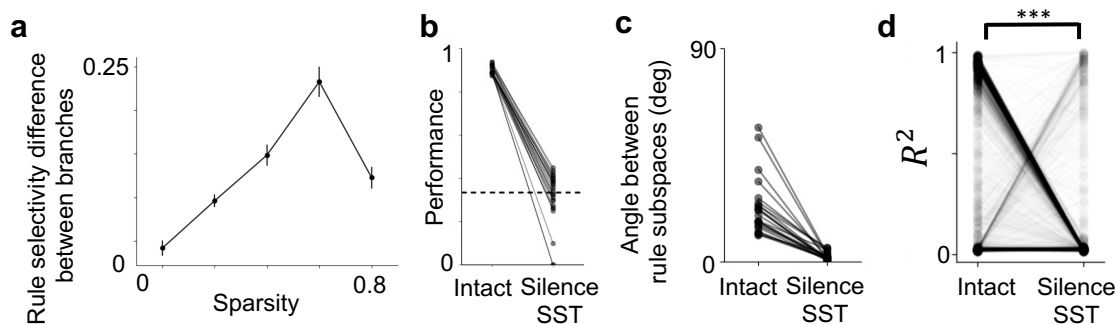


Supplementary Figure 8: The structure in the input and output weights of the sensorimotor module. Data is aggregated across all trained networks with subtractive (a, c) and divisive (b, d) dendritic nonlinearity (c.f. Figure 5b, c). Student's t-test, $p < .001$ for all panels.

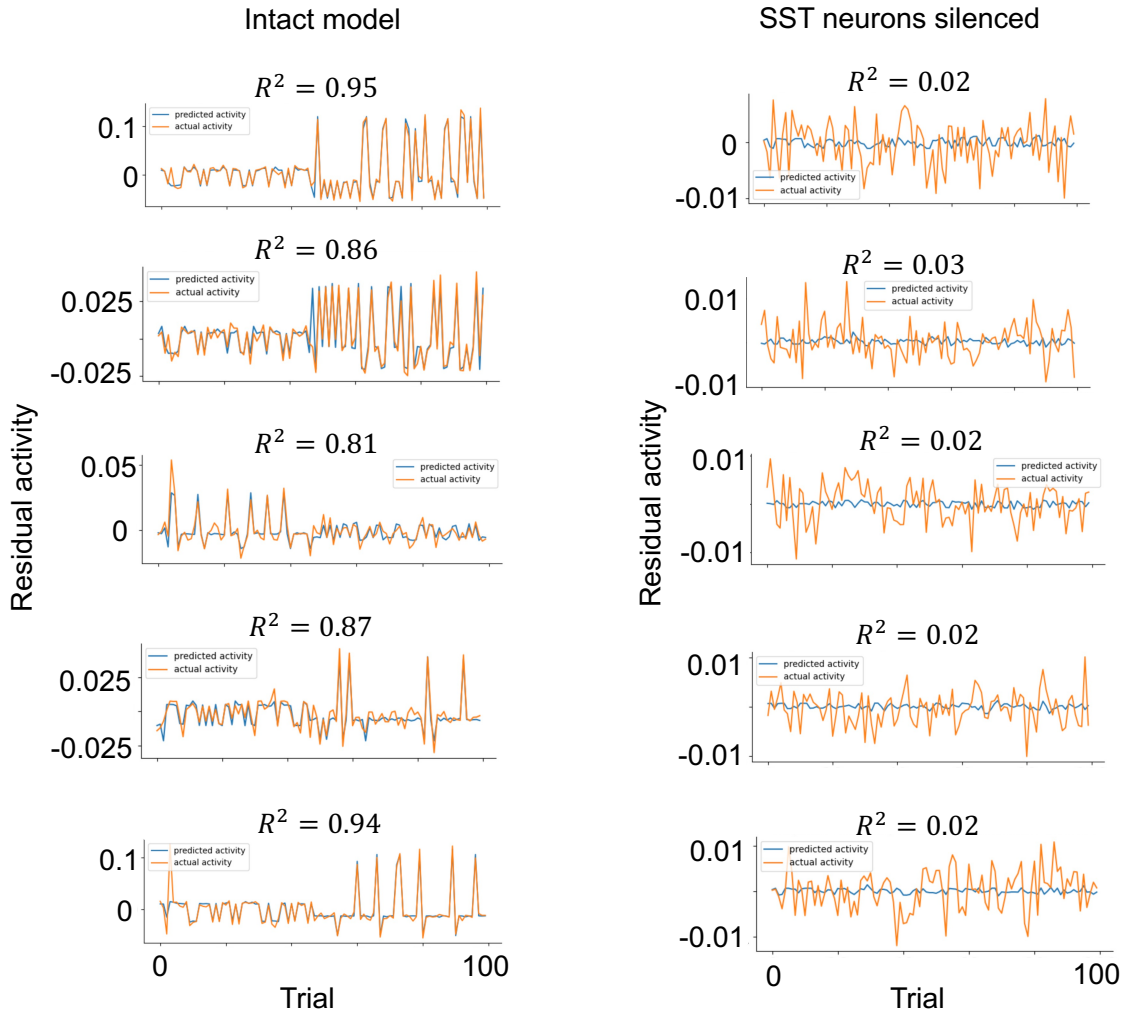
DISSECTING MODULAR RNNs TRAINED TO PERFORM A WCST ANALOG 57



Supplementary Figure 9: a-b. The principal angle between different rule and response subspaces for networks trained with divisive dendritic nonlinearity (c.f. Figure 6c-d). *c-e.* The connectivity bias between different subpopulations of excitatory and PV neurons in the sensorimotor module of networks with divisive dendritic nonlinearity (c.f. Figure 6e-g).



Supplementary Figure 10: a. The relationship between the sparsity of the connection from SST to the dendrites of excitatory neurons in the sensorimotor module, for networks with divisive dendritic nonlinearity (c.f. Figure 7d). *b-d.* The performance (**b**), principal angle between rule subspaces (**c**) and the strength of conjunctive coding (**d**) decreased after silencing SST neurons in the sensorimotor module, for networks trained with divisive dendritic nonlinearity (for **d**, Student's t-test, $p < .001$. Each line represents one neuron. Results are aggregated across networks). c.f. Figure 7e-g.



Supplementary Figure 11: Example neurons in the sensorimotor module showing decreased conjunctive coding of rule and stimulus when SST neurons in the sensorimotor module are silenced. The strength of conjunctive coding is assessed by the R^2 value of a linear regression model where the independent variables are conjunctions of rule and stimulus and the dependent variable is the trial-to-trial residual neural activity unexplained by a model with only rule and stimulus as independent variables (see Methods for details).