

Dynamics of striatal action selection and reinforcement learning

Jack Lindsey¹, Jeffrey E. Markowitz², Sandeep Robert Datta³, and Ashok Litwin-Kumar¹

¹Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY, USA

²Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology and Emory University, Atlanta, GA, USA

³Department of Neurobiology, Harvard Medical School, Boston, MA, USA

Abstract

Spiny projection neurons (SPNs) in dorsal striatum are often proposed as a locus of reinforcement learning in the basal ganglia. Here, we identify and resolve a fundamental inconsistency between striatal reinforcement learning models and known SPN synaptic plasticity rules. Direct-pathway (dSPN) and indirect-pathway (iSPN) neurons, which promote and suppress actions, respectively, exhibit synaptic plasticity that reinforces activity associated with elevated or suppressed dopamine release. We show that iSPN plasticity prevents successful learning, as it reinforces activity patterns associated with negative outcomes. However, this pathological behavior is reversed if functionally opponent dSPNs and iSPNs, which promote and suppress the current behavior, are simultaneously activated by efferent input following action selection. This prediction is supported by striatal recordings and contrasts with prior models of SPN representations. In our model, learning and action selection signals can be multiplexed without interference, enabling learning algorithms beyond those of standard temporal difference models.

Introduction

Numerous studies have proposed that the basal ganglia is a reinforcement learning system (Joel et al., 2002; Niv, 2009; Ito and Doya, 2011). Reinforcement learning algorithms use experienced and predicted rewards to learn to predict the expected future reward associated with an organism's current state and the action to select in order to maximize this reward (Sutton and Barto, 2018). Spiny projection neurons (SPNs) in the striatum are well-positioned to take part in such an algorithm, as they receive diverse contextual information from the cerebral cortex and are involved in both action selection (in dorsal striatum; Packard and Knowlton, 2002; Seo et al., 2012; Balleine et al., 2007) and value prediction (in ventral striatum; Cardinal et al., 2002; Montague et al., 1996; O'Doherty et al., 2004). Moreover, plasticity of SPN input synapses is modulated by midbrain dopamine release (Wickens et al., 1996; Calabresi et al., 2000; Contreras-Vidal and Schultz, 1999). A variety of studies support the view that this dopamine release reflects reward prediction error (Schultz et al., 1997; Montague et al., 1996; Houk and Adams, 1995), which in many reinforcement learning algorithms is the key quantity used to modulate learning (Sutton and Barto, 2018; Niv, 2009).

35 Despite these links, several aspects of striatal physiology are difficult to reconcile with reinforcement
36 learning models. SPNs are classified in two main types – direct-pathway (dSPNs) and indirect-
37 pathway (iSPNs). These two classes of SPNs exert opponent effects on action based on perturbation
38 data (Kravitz et al., 2010; Freeze et al., 2013; Lee and Sabatini, 2021), but also exhibit highly
39 correlated activity (Cui et al., 2013). Moreover, dSPNs and iSPNs express different dopamine
40 receptors (D1-type and D2-type) and thus undergo synaptic plasticity according to different rules.
41 In particular, dSPN inputs are potentiated when coincident pre- and post-synaptic activity is
42 followed by above-baseline dopamine activity, while iSPN inputs are potentiated when coincident
43 pre- and post-synaptic activity is followed by dopamine suppression (Shen et al., 2008; Frank,
44 2005; Iino et al., 2020). Prior studies have attempted to account for these phenomena by proposing
45 that dSPNs learn from positive reinforcement to promote actions, and iSPNs learn from negative
46 reinforcement to suppress actions (Cruz et al., 2022; Collins and Frank, 2014; Jaskir and Frank,
47 2023; Varin et al., 2023). However, we will show that a straightforward implementation of such a
48 model fails to yield a functional reinforcement learning algorithm, as the iSPN learning rule assigns
49 blame for negative outcomes to the wrong actions.

50 In this work, we begin by rectifying this inconsistency between standard reinforcement learning
51 models of the striatum and known SPN plasticity rules. The iSPN learning rule reported in the
52 literature reinforces patterns of iSPN activity that are associated with dopamine suppression, in-
53 creasing the likelihood of repeating decisions that previously led to negative outcomes. We show
54 that this pathological behavior is reversed if opponent dSPNs and iSPNs receive correlated efferent
55 input encoding the animal’s present action. This model provides an explanation for the apparent
56 paradox that the activities of dSPNs and iSPNs are positively correlated despite their opponent
57 causal functions (Cui et al., 2013). Importantly, our model predicts coactivity following action se-
58 lection of dSPNs and iSPNs that are responsible for regulating the same behavior (promoting and
59 suppressing it, respectively). This somewhat counterintuitive prediction contrasts with prior pro-
60 posals that dSPNs that promote an action are coactive with iSPNs that suppress different actions
61 (Mink, 1996; Redgrave et al., 1999). We find support for this prediction in experimental recordings
62 of dSPNs and iSPNs during spontaneous behavior.

63 Next, we show that the sign difference between dSPN and iSPN plasticity rules enables more
64 sophisticated learning algorithms than can be achieved in models with a single plasticity rule.
65 Specifically, heterogeneity of SPN plasticity rules allows action selection signals to be multiplexed
66 with feedforward SPN activity without interference. This enables the striatum to implement so-
67 called *off-policy* reinforcement learning algorithms, in which the cortico-striatal pathway learns
68 from the the outcomes of actions that are driven by other neural pathways. Off-policy algorithms
69 are commonly used in state-of-the-art machine learning algorithms, as they dramatically improve
70 learning efficiency by facilitating learning from expert demonstrations, mixture-of-experts models,
71 and replayed experiences (Arulkumaran et al., 2017). Following the implications of this model
72 further, we show that off-policy algorithms require a dopaminergic signal in dorsal striatum that
73 combines classic state-based reward prediction error with a form of action prediction error. We
74 confirm a key signature of this prediction in recent dopamine data collected from dorsolateral
75 striatum during spontaneous behavior.

76 Results

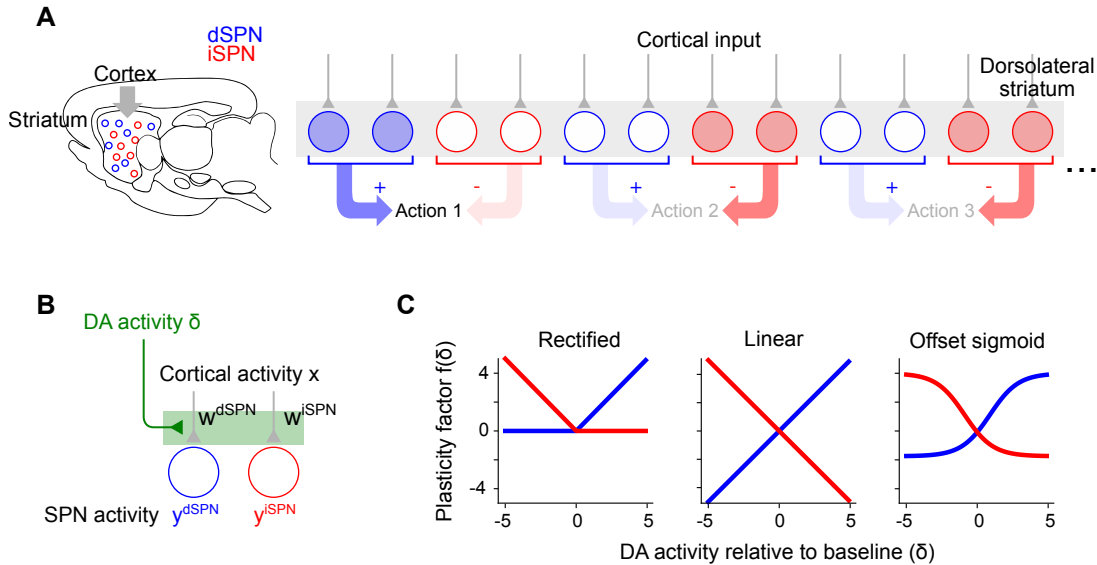


Figure 1: Corticostriatal action selection circuits and plasticity rules. **A**. Left, diagram of cortical inputs to striatal populations. Right, illustration of action selection architecture. Populations of dSPNs (blue) and iSPNs (red) in DLS are responsible for promoting and suppressing specific actions, respectively. Active neurons (shaded circles) illustrate a pattern of activity consistent with typical models of striatal action selection, in which dSPNs that promote a chosen action and iSPNs that suppress other actions are active. **B**. Illustration of three-factor plasticity rules at SPN input synapses, in which adjustments to corticostriatal synaptic weights depend on presynaptic cortical activity, SPN activity, and dopamine release. **C**. Illustration of different models of the dopamine-dependent factor $f(\delta)$ in dSPN (blue) and iSPN (red) plasticity rules.

In line with previous experimental (Wickens et al., 1996; Calabresi et al., 2000; Contreras-Vidal and Schultz, 1999) and modeling (Sutton and Barto, 2018; Niv, 2009) studies, we model plasticity of corticostriatal synapses using a three-factor learning rule, dependent on coincident presynaptic activity, postsynaptic activity, and dopamine release (Fig 1A,B). Concretely, we model plasticity of the weight w of a synapse from a cortical neuron with activity x onto a dSPN or iSPN with activity y as

$$w^{\text{dSPN}} = f^{\text{dSPN}}(\delta) \cdot y^{\text{dSPN}} \cdot x, \quad (1)$$

$$w^{\text{iSPN}} = f^{\text{iSPN}}(\delta) \cdot y^{\text{iSPN}} \cdot x, \quad (2)$$

77 where δ represents dopamine release relative to baseline, and the functions $f^{\text{dSPN}}(\delta)$ and $f^{\text{iSPN}}(\delta)$
 78 model the dependence of the two plasticity rules on dopamine concentration.

79 For dSPNs, the propensity of input synapses to potentiate increases with increasing dopamine
 80 concentration, while for iSPNs the opposite is true. This observation is corroborated by converging
 81 evidence from observations of dendritic spine volume, intracellular PKA measurements, and spike-
 82 timing dependent plasticity protocols (Shen et al., 2008; Gurney et al., 2015; Iino et al., 2020;
 83 Lee et al., 2021). For the three-factor plasticity rule above, these findings imply that f^{dSPN} is an
 84 increasing function of δ while f^{iSPN} is a decreasing function. Prior modeling studies have proposed

85 specific plasticity rules that correspond to different choices of f^{dSPN} and f^{iSPN} , some examples of
86 which are shown in Fig. 1C.

87 **iSPN plasticity rule impedes successful reinforcement learning**

88 Prior work has proposed that dSPNs activate when actions are performed and iSPNs activate when
89 actions are suppressed (Fig. 1A). When an animal selects among multiple actions, subpopulations
90 of dSPNs are thought to promote the selected action, while other subpopulations of iSPNs inhibit
91 the unchosen actions (Mink, 1996; Redgrave et al., 1999). We refer to this general description as
92 the “canonical action selection model” of SPN activity and show that this model, when combined
93 with the plasticity rules above, fails to produce a functional reinforcement learning algorithm.
94 This failure is specifically due to the iSPN plasticity rule. Later, we also show that the SPN
95 representation predicted by the canonical action selection model is inconsistent with recordings of
96 identified dSPNs and iSPNs. We begin by analyzing a toy model of a “go/no-go” task in which
97 the “go” action is rewarded. In the model, the probability of selecting the “go” action is increased
98 when a dSPN is active and decreased when an iSPN is active (Fig. 2A). After an action is taken,
99 dopamine activity reports the reward prediction error, increasing when reward is obtained and
100 decreasing when it is not.

101 It is easy to see that the dSPN plasticity rule in Eq. (1) is consistent with successful reinforcement
102 learning (Fig. 2A). Suppose a “go” action is selected due to increased activity of the dSPN, leading
103 to reward (Fig. 2A, center). The resulting dopamine increase potentiates inputs to the dSPN from
104 cortical neurons that are active during the task. As a result, the “go” action is more likely to be
105 selected in the future.

106 At first glance, it may seem that a similar logic would apply to iSPNs, since their suppressive effect
107 on behavior and reversed dependence on dopamine concentration are both opposite to dSPNs.
108 However, a more careful examination reveals that the iSPN plasticity rule in Eq. (2) does not
109 promote successful learning. If the “no-go” action is selected and reward is not obtained (Fig. 2B,
110 center), the resulting dopamine decrease and iSPN plasticity rule increases the tendency of the iSPN
111 to inhibit the “go” action, making future reward even less likely (Fig. 2B, right). More generally,
112 the model demonstrates that, while the plasticity rule of Eq. (1) correctly reinforces dSPN activity
113 patterns that lead to positive outcomes, it incorrectly reinforces iSPN activity patterns that lead
114 to negative outcomes. The function of iSPNs in inhibiting action does not change the fact that
115 such reinforcement is undesirable.

116 These pathological learning dynamics are also evident in tasks that require selecting between mul-
117 tiple actions. In the canonical action selection model, dSPNs promoting a selected action and
118 iSPNs inhibiting unselected actions are active. If a negative outcome is encountered leading to a
119 dopamine decrease, Eq. (1) predicts that inputs to iSPNs corresponding to unselected actions are
120 strengthened (LTP in Fig. 2C, left). This makes the action that led to the negative outcome more
121 rather than less likely to be taken when the same cortical inputs are active in the future (Fig. 2C,
122 right). We note that, depending on the learning rule (Fig. 1C), inputs to dSPNs that promote the
123 selected action may be weakened (LTD in Fig. 2C, left), which correctly disincentivizes the action
124 that led to a negative outcome. However, this dSPN effect competes with the pathological behavior
125 of the iSPNs and is often unable to overcome it.

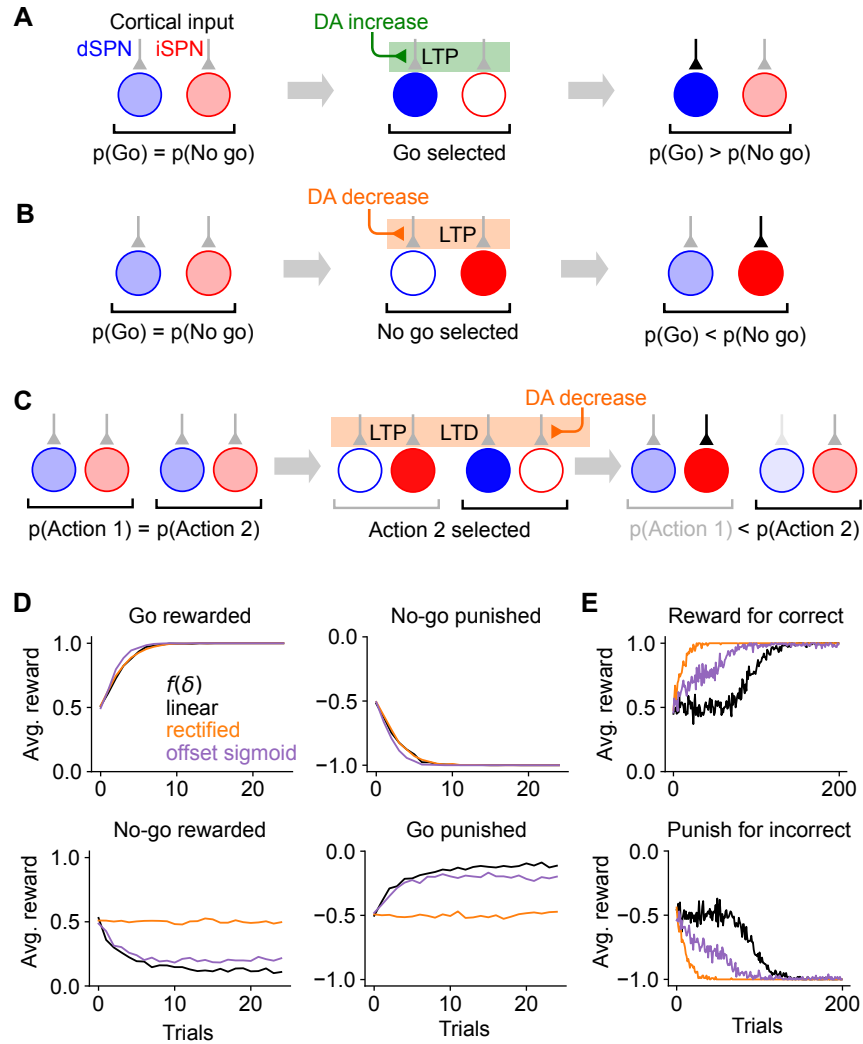


Figure 2: Consequences of the canonical action selection model of SPN activity. **A.** Example in which dSPN plasticity produces correct learning behavior in a go/no-go task. Left: cortical inputs to the dSPN and iSPN are equal prior to learning. Middle: the “go” response is selected, corresponding to elevated dSPN activity. In this example, the “go” response is rewarded, leading to elevated DA activity and thus potentiation of the dSPN input synapse. Right: in a subsequent trial, cortical input to the dSPN is stronger, increasing the likelihood of selecting the “go” response. **B.** Example in which iSPN plasticity produces incorrect learning behavior in a go/no-go task. Left: same as panel B. Middle: the “no go” response is selected, corresponding to elevated iSPN activity. In this example, the “no-go” response is punished, leading to decreased DA activity and thus potentiation of the iSPN input synapse. Right: in a subsequent trial, cortical input to the iSPN is stronger, decreasing the likelihood of selecting the “go” response. **C.** Example in which iSPN plasticity produces incorrect learning behavior in an action selection task. Left: cortical inputs to all SPNs are equal prior to learning. Middle: action 2 is selected, corresponding to elevated activity in the dSPN that promotes action 2 and the iSPN that suppresses action 1. In this example, action 2 is punished, leading to decreased DA activity. The input synapse to the action 2-promoting dSPN is (depending on the learning rule, see Fig. 1) depressed, and the input to the action 1-suppressing iSPN is potentiated. Right: On a subsequent trial, input to the action 1-suppressing iSPN is stronger, decreasing the probability of selecting action 1 rather than action 2. Note that the dSPN input corresponding to action 2 is (potentially) weakened, which correctly decreases the probability of selecting action 2, but this effect is not sufficient to overcome the strengthened action 1 iSPN activity. **D.** Performance of a simulated striatal reinforcement learning system in go/no-go tasks with different reward contingencies. **E.** Same as D, but for action selection tasks with two cortical input states, two available actions, and one correct action per state, under different reward protocols.

126 We also note that, if dopamine increases lead to depression of iSPN inputs (Fig. 1C, center, right),
 127 positive outcomes will lead to actions that were correctly being inhibited by iSPNs to be less
 128 inhibited in the future. Thus, both positive and negative outcomes may cause incorrect iSPN
 129 learning. Some sources suggest that while dopamine suppression increases D2 receptor activation,
 130 dopamine increase has little effect on D2 receptors (Dreyer et al., 2010), corresponding to the
 131 rectified model of $f(\delta)$ (Fig. 1C, left). In this case, pathological iSPN plasticity behavior still
 132 manifests when dopamine activity is suppressed (as in the examples of Fig. 2B,C).

133 We simulated learning of the go/no-go task with the three-factor plasticity rules above, with
 134 dopamine activity modeled as reward prediction error obtained using a temporal difference learning
 135 rule. As expected, the system learns the wrong behavior when performance feedback is provided
 136 on no-go trials, and thus iSPN plasticity is the main driver of learning (Fig. 2D). We also simulated
 137 a two-alternative forced choice task in which there are two cues (corresponding to different cortical
 138 input patterns), each with a corresponding target action. When performance feedback consists
 139 of rewards for correct actions, the system learns the task, as dSPNs primarily drive the learning.
 140 However, when instead performance feedback consists of giving punishments for incorrect actions,
 141 the system does not learn the task, as iSPNs primarily drive the learning (Fig. 2E). We note that,
 142 in principle, this problem could be avoided if the learning rate of iSPNs were very small compared
 143 to that of dSPNs, ensuring that reinforcement learning is always primarily driven by the dSPN
 144 pathway (leaving iSPNs to potentially perform a different function). However, this alternative
 145 would be inconsistent with prior studies indicating a significant role for the indirect pathway in
 146 reinforcement learning (Peak et al., 2020; Lee and Sabatini, 2021). The model we introduce below
 147 makes use of contributions to learning from both pathways.

148 **Efferent activity in SPNs enables successful reinforcement learning**

We have shown that the canonical action selection model, when paired with Eq. (1), produces incorrect learning. What pattern of SPN activity would produce correct learning? In the model, the probability of selecting an action is determined by the “difference mode” $y^{\text{dSPN}} - y^{\text{iSPN}}$, where y^{dSPN} and y^{iSPN} are the activities of dSPN and iSPN neurons associated with that action. We analyzed how the plasticity rule of Eq. (1) determines changes to this difference mode. In the simplest case in which the SPN firing rate is a linear function of cortical input (that is, $y^{\text{d/iSPN}} = \mathbf{w}^{\text{d/iSPN}} \cdot \mathbf{x}$) and plasticity’s dependence on dopamine concentration is also linear (that is, $f^{\text{d/iSPN}}(\delta) \propto \pm\delta$; Fig. 1C, center), the change in the probability of selecting an action due to learning is

$$\begin{aligned}
 \Delta(y^{\text{dSPN}} - y^{\text{iSPN}}) &= \Delta\mathbf{w}^{\text{dSPN}} \cdot \mathbf{x} - \Delta\mathbf{w}^{\text{iSPN}} \cdot \mathbf{x} \\
 &\propto \delta y^{\text{dSPN}}(\mathbf{x} \cdot \mathbf{x}) - (-\delta)y^{\text{iSPN}}(\mathbf{x} \cdot \mathbf{x}) \\
 &\propto \delta(y^{\text{dSPN}} + y^{\text{iSPN}}).
 \end{aligned}
 \tag{3}$$

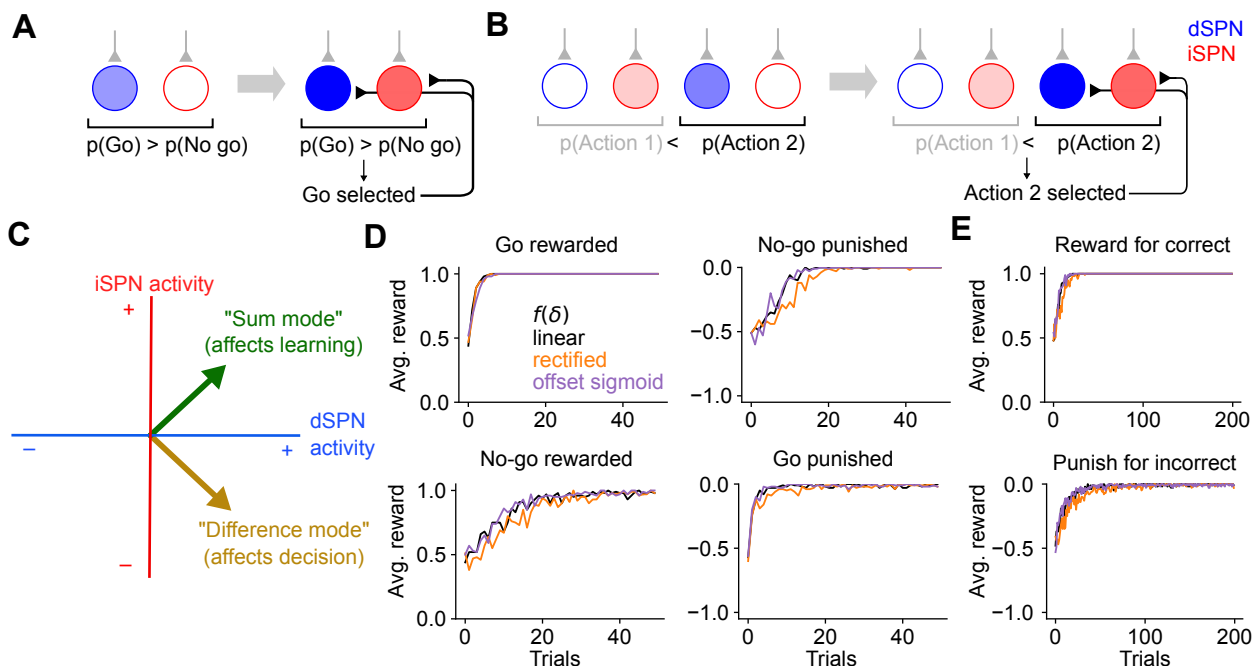


Figure 3: The efference model of SPN activity. **A.** Illustration of the efference model in a go/no-go task. Left: feedforward SPN activity driven by cortical inputs. Right: once the “go” response is selected, the dSPN and iSPN are both excited by efferent input, which is combined with their original input. As a result, both the dSPN and iSPN are more active than prior to action selection, but the dSPN is still more active than the iSPN. **B.** Illustration of the efference model in an action selection task. Left: feedforward SPN activity driven by cortical inputs. Right: once action 2 is selected, efferent inputs excite the dSPN and iSPN responsible for promoting and suppressing action 2. Efferent activity is combined with feedforward activity, such that the action 2-associated dSPNs and iSPNs are both more active than the action 1 dSPNs and iSPNs, but the relative dSPN and iSPN activity for each action remains unchanged. **C.** The activity levels of the dSPN and iSPN populations that promote and suppress a given action can be plotted in a two-dimensional space. The difference mode influences the probability of taking that action, while activity in the sum mode drives future changes to activity in the difference mode via plasticity. Efferent activity excites the sum mode. **D.** Performance of a striatal RL system using the efference model on the tasks of Fig. 2D. **E.** Performance of a striatal RL system using the efference model on the tasks of Fig. 2E.

149 Changes to the “difference mode” $y^{\text{dSPN}} - y^{\text{iSPN}}$ are therefore driven by the “sum mode” $y^{\text{dSPN}} +$
 150 y^{iSPN} . This implies that the activity pattern that leads to correct learning about an action’s outcome
 151 is different from the activity pattern that selects the action. To promote or inhibit, respectively, an
 152 action that leads to a dopamine increase or decrease, this analysis predicts that both dSPNs that
 153 promote and iSPNs that inhibit the action should be co-active. A more general argument applies
 154 for other learning rules and firing rate nonlinearities: as long as $y^{\text{d/iSPN}}$ is an increasing function
 155 of total input current, $f^{\text{dSPN}}(\delta)$ has positive slope, and $f^{\text{iSPN}}(\delta)$ has negative slope, changes in
 156 difference mode activity will be positively correlated with sum mode activity (see Supplemental
 157 Information).

158 The key insight of the above argument is that the pattern of SPN activity needed for learning
 159 involves simultaneous excitation of dSPNs that promote the current behavior and iSPNs that
 160 inhibit it. This differs from the pattern of activity needed to drive selection of that behavior
 161 in the first place. We therefore propose a model in which SPN activity contains a substantial
 162 *efferent* component that follows action selection and promotes learning, but has no causal impact
 163 on behavior. In the model, feedforward cortico-striatal inputs initially produce SPN activity whose

164 difference mode causally influences action selection, consistent with the canonical model (Fig. 3A,B,
165 left). When an action is performed, both dSPNs and iSPNs responsible for promoting or inhibiting
166 that action receive efferent excitatory input, producing sum-mode activity. Following this step, SPN
167 activity reflects both contributions (Fig. 3A,B, right). Unlike the canonical action selection model
168 (Fig. 1A), this model thus predicts an SPN representation in which, after an action is selected,
169 the most highly active neurons are those responsible for regulating that behavior and not other
170 behaviors.

171 In SPN activity space, the sum and difference modes are orthogonal to one another. This orthog-
172 onality has two consequences. First, it implies that encoding the action in the difference mode (as
173 in the canonical action selection model) produces synaptic weight changes that do not promote
174 learning, consistent with the competing effects of dSPN and iSPN plasticity that we previously
175 described. Second, it implies that adding efferent activity along the sum mode, which produces
176 correct learning, has no effect on action selection. The model thus provides a solution to the
177 problem of interference between “forward pass” (action selection) and “backward pass” (learning)
178 activity, a common issue in models of biologically plausible learning algorithms (see Discussion).

179 In simulations, we confirm that unlike the canonical action selection model, this efference model
180 solves go/no-go (Fig. 3D) and action selection (Fig. 3E) tasks regardless of the reward protocol.
181 Although the derivation above assumes linear SPN responses and linear dependences of plasticity
182 on dopamine concentration, our model enables successful learning even using a nonlinear model
183 of SPN responses and a variety of plasticity rules (Fig. 3D,E; see Supplemental Information for a
184 derivation that explains this general success). Finally, we also confirmed that our results apply to
185 cases in which actions are associated with distributed modes of dSPN and iSPN activity, and with
186 a larger action space (Supp. Fig. 1).

187 Efference model predicts properties of SPN activity

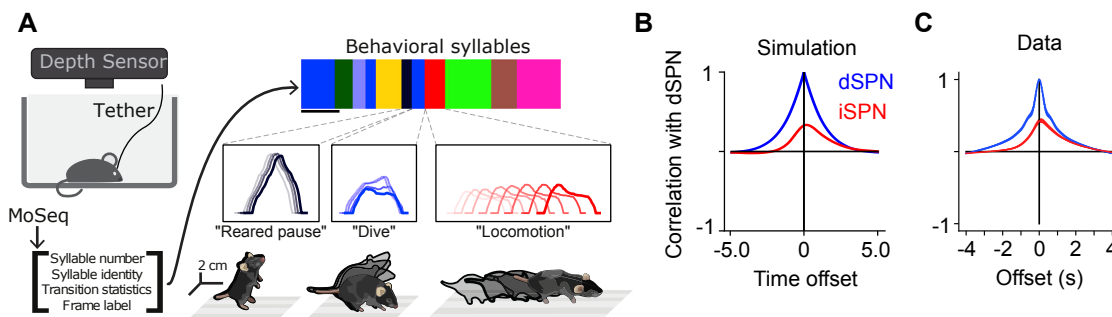


Figure 4: Comparisons of model predictions about bulk dSPN and iSPN activity to experimental data. **A.** Schematic of experimental setup, taken from Markowitz et al. (2018). Neural activity and kinematics of spontaneously behaving mice are recorded, and behavior is segmented into stereotyped “behavioral syllables” using the MoSeq pipeline. **B.** In simulation of efference model with random feedforward cortical inputs, cross-correlation of total dSPN and iSPN activity. **C.** Cross-correlation between fiber photometry recordings of bulk dSPN and iSPN activity in freely behaving mice, using the data from Markowitz et al. (2018). Line thickness indicates standard error of the mean.

188 Thus far, we have provided theoretical arguments and model simulations that suggest that simul-
189 taneous efferent input to opponent dSPNs and iSPNs is necessary for reinforcement learning, given

190 known plasticity rules. We next sought to test this prediction in neural data. We used data from a
191 recent study which recorded bulk and cellular dSPN and iSPN activity in spontaneously behaving
192 mice (Fig. 4A; Markowitz et al., 2018). As no explicit rewards or task structure were provided
193 during recording sessions, we adopted a modeling approach that makes minimal assumptions about
194 the inputs to SPNs besides the core prediction of efferent activity. Specifically, we used a network
195 model in which (1) populations of dSPNs and iSPNs promote or suppress different actions, (2) the
196 feedforward inputs to all SPNs are random, (3) actions are sampled with log-likelihoods scaling
197 according to the associated dSPN and iSPN difference mode, and (4) efferent activity excites the
198 sum mode corresponding to the chosen action.

199 In this model, difference mode dSPN and iSPN activity drives behaviors, and those behaviors cause
200 efferent activation of the corresponding sum mode. As a result, on average, dSPN activity tends to
201 lead to increased future iSPN activity, while iSPN activity leads to decreased future dSPN activity.
202 Consequently, the temporal cross-correlation between total dSPN activity and iSPN activity is
203 asymmetric, with present dSPN activity correlating more strongly with future iSPN activity than
204 with past iSPN activity (Fig. 4B). Such asymmetry is not predicted by the canonical action selection
205 model, or models that assume dSPNs and iSPNs are co-active. Computing the temporal cross-
206 correlation in the bulk two-color photometry recordings of dSPN and iSPN activity, we find a very
207 similar skewed relationship in the data (Fig. 4C). We confirmed this result is not an artifact of the
208 use of different indicators for dSPN and iSPN activity by repeating the analysis on data from mice
209 where the indicators were reversed and finding the same result (Supp. Fig. 2).

210 Our model makes even stronger predictions about SPN population activity and its relationship to
211 action selection. First, it predicts that both dSPNs and iSPNs exhibit similar selectivity in their
212 tuning to actions. This contrasts with implementations of the canonical action selection model in
213 which iSPNs are active whenever their associated action is not being performed and thus are more
214 broadly tuned than dSPNs (Fig. 1A). Second, it also predicts that efferent activity excites dSPNs
215 that promote the currently performed action and iSPNs that suppress the currently performed
216 action. As a result, dSPNs whose activity increases during the performance of a given action
217 should tend to be above baseline shortly prior to the performance of that action. By contrast,
218 iSPNs whose activity increases during an action should tend to be below baseline during the same
219 time interval (Fig. 5A, left). Moreover, this effect should be action-specific: the dSPNs and iSPNs
220 whose activity increases during a given action should display negligible average fluctuations around
221 the onset of other actions (Fig. 5A, right). These predictions can also be reinterpreted in terms of
222 the sum and difference modes. The difference mode activity associated with an action is elevated
223 prior to selection of that action, while the sum mode activity is excited following action selection
224 (Fig. 5B).

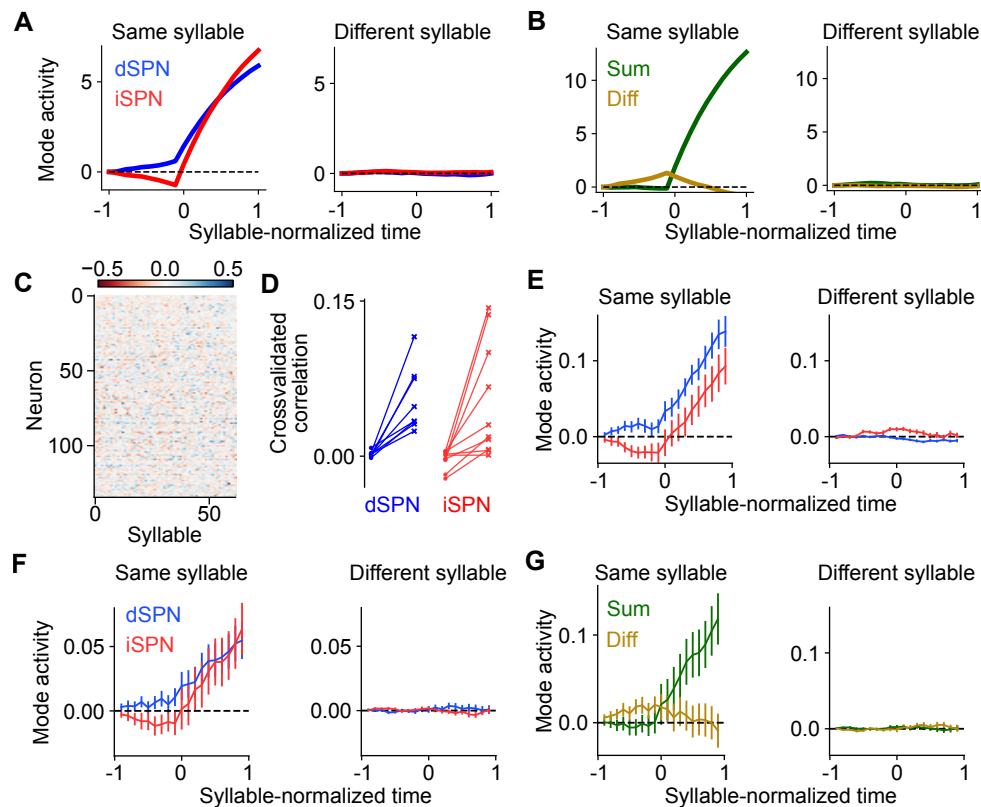


Figure 5: Comparisons of model predictions about action-tuned SPN subpopulations to experimental data. **A.** Activity of dSPNs (blue) and iSPNs (red) around the onset of their associated action (left) or other actions (right) in the simulation from Fig. 4. **B.** Same information as A, but plotting activity of the sum (dSPN + iSPN) and difference (dSPN - iSPN) modes. **C.** For an example experimental session, dSPN activity modes associated with each of the behavioral syllables, in z-scored firing rate units. **D.** Correlation between identified dSPN and iSPN activity modes in two random subsamples of the data, for shuffled (left, circles) and real (right, x's) data. **E.** Projection of dSPN (blue) and iSPN (red) activity onto the syllable-associated modes identified in panel C, around the onset of the associated syllable (left panel) or other syllables (right panel) averaged across all syllables. Error bars indicate standard error of the mean across syllables. **F.** Same as panel E, restricting the analysis to mice in which dSPNs and iSPNs were simultaneously recorded. **G.** Same data as panel F, but plotting activity of the sum (dSPN + iSPN) and difference (dSPN - iSPN) modes.

225 To test these hypotheses, we used calcium imaging data collected during spontaneous mouse behavior (Markowitz et al., 2018). The behavior of the mice was segmented into consistent, stereotyped
 226 kinematic motifs referred to as “behavioral syllables,” as in previous studies (Fig. 4A). We regard
 227 these behavioral syllables as the analogs of actions in our model. First, we examined the tuning
 228 of dSPNs and iSPNs to different actions and found that, broadly consistent with what our model
 229 predicts, both subpopulations exhibit similar selectivities (Supp. Fig. 3). Next, to test our predic-
 230 tions about dynamics before and after action selection (Fig. 5A,B), we identified, for each syllable,
 231 dSPN and iSPN population activity vectors (“modes”) that increased the most during performance
 232 of that syllable (Fig. 5C). We confirmed that these modes are meaningful by checking that modes
 233 identified using two disjoint subsets of the data are correlated (Fig. 5D). We then plotted the activ-
 234 ity of these modes around the time of onset of the corresponding syllable, and averaged the result
 235 across the choice of syllables (Fig. 5E). The result displays remarkable agreement with the model
 236 prediction in Fig. 5A.
 237

238 The majority of the above data consisted of recordings of either dSPNs or iSPNs from a given
 239 mouse. However, in a small subset ($n=4$) of mice, dSPNs and iSPNs were simultaneously recorded
 240 and identified. We repeated the analysis above on these sessions, and found the same qualitative
 241 results (Fig. 5F). The simultaneous recordings further allowed us to visualize the sum and difference
 242 mode activity (Fig. 5G), which also agrees with the predictions of our model (Fig. 5B).

243 Efference model enables off-policy reinforcement learning

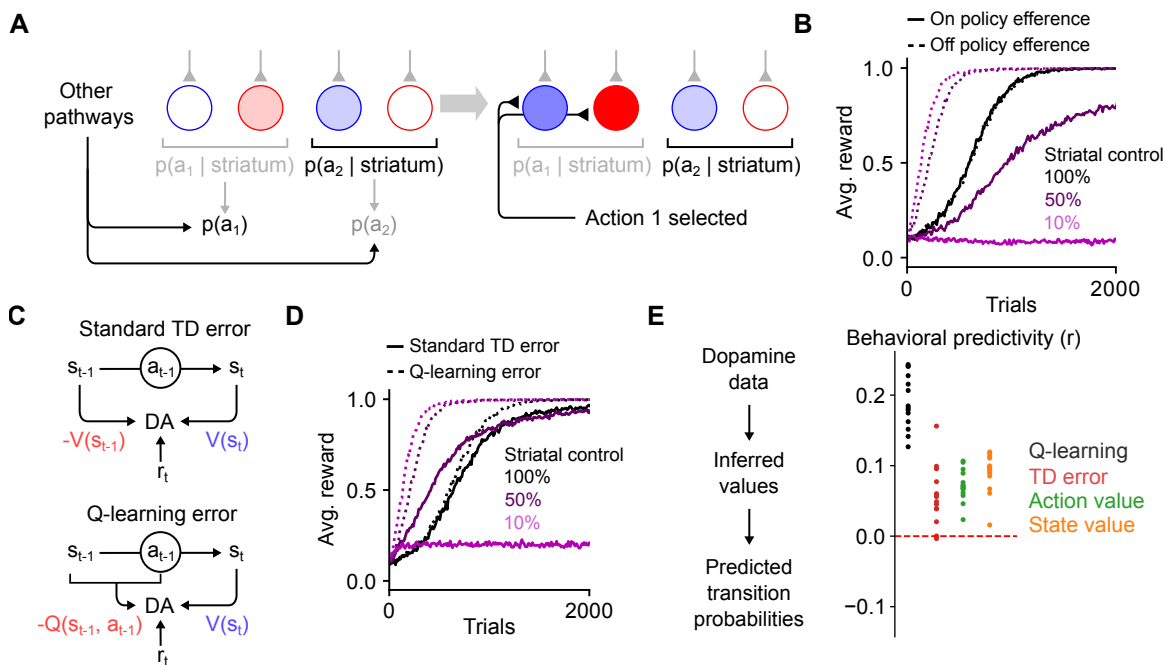


Figure 6: The efference model enables off-policy reinforcement learning. **A.** Illustration of the efference model when the striatum shares control of behavior with other pathways. In this example, striatal activity biases the action selection toward choosing action 2, but other neural pathways override the striatum and cause action 1 to be selected instead (left). Following action selection, efferent activity excites the dSPN and iSPN associated with action 1. However, the outputs of the striatal population remain unchanged. **B.** Performance of RL models in a simulated action selection task (10 cortical states, 10 available actions, in each state one of the actions results in a reward of 1 and the others result in zero reward). Control is shared between the striatal RL circuit and another pathway that biases action selection toward the correct action. Different lines indicate different strength of striatal control relative to the strength of the other pathway. Line style (dashed or solid) indicates the efference model: off-policy efference excites SPNs associated with the selected action, while on-policy efference excites SPNs associated with the action most favored by the striatum. **C.** Schematic of different reinforcement learning models of dopamine activity. The standard TD error models predicts that dopamine activity is sensitive to reward, the predicted value of the current state, and the predicted value of the previous state. The Q-learning error model predicts sensitivity to reward, the predicted value of the current state, and the predicted value of the previous state-action pair. **D.** In the task of panel B using the off-policy efference model, comparison between different models of dopamine activity as striatal control is varied (the Q-learning error model was used in panel B). **E.** Correlation between predicted and actual syllable-to-syllable transition matrix. Predictions were made according to different models of the relationship between dopamine activity and behavior, using observed average dopamine activity associated with syllable transitions in the data of Markowitz et al. (2023). Each dot indicates a different experimental session.

244 Prior studies have argued for the importance of motor efference copies during basal ganglia learn-
245 ing, in particular when action selection is influenced by other brain regions (Fee, 2014; Lindsey
246 and Litwin-Kumar, 2022). Indeed, areas such as the motor cortex and cerebellum drive behavior
247 independent of the basal ganglia (Exner et al., 2002; Wildgruber et al., 2001; Ashby et al., 2010;
248 Silveri, 2021; Bostan and Strick, 2018). Actions taken by an animal may therefore at times differ
249 from those most likely to be selected by striatal outputs (Fig. 6A), and it may be desirable for
250 cortico-striatal synapses to learn about the consequences of these actions.

251 In the reinforcement learning literature, this kind of learning is known as an “off-policy” algorithm,
252 as the reinforcement learning system (in our model, the striatum) learns from actions that follow
253 a different policy than its own. Off-policy learning has been observed experimentally, for instance
254 in the consolidation of cortically driven behaviors into subcortical circuits including dorsolateral
255 striatum (Kawai et al., 2015; Hwang et al., 2019; Mizes et al., 2023). Such learning requires efferent
256 activity in SPNs that reflects the actions being performed, rather than the action that would be
257 performed based on the striatum’s influence alone.

258 We modeled this scenario by assuming that action selection is driven by weighted contributions from
259 both the striatum and other motor pathways and that the ultimately selected action drives efferent
260 activity (Fig. 6A; see Methods). We found that when action selection is not fully determined by the
261 striatum, such efferent activity is critical for successful learning (Fig. 6B). Notably, in our model,
262 efferent activity has no effect on striatal action selection, due to the orthogonality of the sum and
263 difference modes (Fig. 3C). In a hypothetical alternative model in which the iSPN plasticity rule
264 is the same as that of dSPNs, the efferent activity needed for learning is not orthogonal to the
265 output of the striatum, impairing off-policy learning (Supp. Fig. 4). Thus, efferent excitation of
266 opponent dSPNs/iSPNs is necessary both to implement correct learning updates given dSPN and
267 iSPN plasticity rules, and to enable off-policy reinforcement learning.

268 **Off-policy reinforcement learning predicts relationship between dopamine activ-** 269 **ity and behavior**

We next asked whether other properties of striatal dynamics are consistent with off-policy reinforcement learning. We focused on the dynamics of dopamine release, as off-policy learning makes specific predictions about this signal. Standard temporal difference (TD) learning models of dopamine activity (Fig. 6C, top) determine the expected future reward (or “value”) $V(s)$ associated with each state s using the following algorithm:

$$\delta_t = r_t + V(s_t) - V(s_{t-1}) \quad (4)$$

$$V(s_t) \leftarrow V(s_t) + \alpha \delta_t, \quad (5)$$

270 where s_t and s_{t-1} indicate current and previous states, r_t indicates the currently received reward,
271 α is a learning rate factor, and δ_t is the TD error thought to be reflected in phasic dopamine
272 responses. These dopaminergic responses can be used as the learning signal for a updating action
273 selection in dorsal striatum (Eq. 1, 2), an arrangement commonly referred to as an “actor-critic”
274 architecture (Niv, 2009).

TD learning is an on-policy algorithm, in that the values $V(s)$ associated with each state are calculated under the assumption that the system’s future actions will be similar to those taken

during learning. Hence, TD learning is poorly suited to training an action selection policy in the striatum in situations where the striatum does not fully control behavior, as the values $V(s)$ will not reflect the expected future reward associated with a state if the striatum were to dictate behavior on its own. Off-policy algorithms such as Q-learning solve this issue by learning an action-dependent value function $Q(s, a)$, which indicates the expected reward associated with taking action a in action s (Fig. 6C, bottom), via the following algorithm:

$$\delta_t = r_t + V(s_t) - Q(s_{t-1}, a_{t-1}) \quad (6)$$

$$V(s) = \max_a Q(s, a). \quad (7)$$

275 This algorithm predicts that the dopamine response δ_t is action-dependent. The significance of on-
276 policy vs. off-policy learning algorithms can be demonstrated in simulations of operant conditioning
277 tasks in which control of action selection is shared between the striatum and another “tutor”
278 pathway that biases responses toward the correct action. When the striatal contribution to decision-
279 making is weak, it is unable to learn the appropriate response when dopamine activity is modeled
280 as a TD error (Fig. 6D). On the other hand, a Q-learning model of dopamine activity enables
281 efficient striatal learning even when control is shared with another pathway.

282 For the spontaneous behavior paradigm we analyzed previously (Fig. 4A), Q-learning but not
283 TD learning predicts sensitivity of dopamine responses to the likelihood of the previous syllable-
284 to-syllable transition. Using recordings of dopamine activity in the dorsolateral striatum in this
285 paradigm (Markowitz et al., 2023), we tested whether a Q-learning model could predict the relation-
286 ship between dopamine activity and behavioral statistics, comparing it to TD learning and other
287 alternatives (see Supplemental Information). The Q-learning model matches the data significantly
288 better than alternatives (Fig. 6E), providing support for a model of dorsal striatum as an off-policy
289 reinforcement learning system.

290 Discussion

291 We have presented a model of reinforcement learning in the dorsal striatum in which efferent ac-
292 tivity excites dSPNs and iSPNs that promote and suppress, respectively, the currently selected
293 action. Thus, following action selection, iSPN activity counterintuitively represents the action that
294 is inhibited by the currently active iSPN population. This behavior contrasts with previous pro-
295 posals, in which iSPN activity reflects actions being inhibited. This model produces updates to
296 corticostriatal synaptic weights given the known opposite-sign plasticity rules in dSPNs and iSPNs
297 that correctly implement a form of reinforcement learning (Fig. 3), which in the absence of such
298 efferent activity produce incorrect weight updates (Fig. 2). The model makes several novel pre-
299 dictions about SPN activity which we confirmed in experimental data (Figs. 4, 5). It also enables
300 multiplexing of action selection signals and learning signals without interference. This facilitates
301 more sophisticated learning algorithms such as off-policy reinforcement learning, which allows the
302 striatum to learn from actions that were driven by other neural circuits. Off-policy reinforcement
303 learning requires dopamine to signal action-sensitive reward predictions errors, which agrees better
304 with experimental recordings of striatal dopamine activity than alternative models (Fig. 6).

305 Other models of efferent inputs to the striatum

306 Prior work has pointed out the need for efference copies of decisions to be represented in the
307 striatum, particularly for actions driven by other circuits (Fee, 2014). Frank (2005) propose a model
308 in which premotor cortex outputs collateral signals to the striatum that represent the actions under
309 consideration, with the striatum potentially biasing the decision based on prior learning. Through
310 bidirectional feedback (premotor cortex projecting to striatum, and striatum projecting to premotor
311 cortex indirectly through the thalamus) a decision is collectively made by the combined circuit, and
312 the selected action is represented in striatal activity, facilitating learning about the outcome of the
313 action. While similar to our proposal in some ways, this model implicitly assumes that the striatal
314 activity necessary for decision-making is also what is needed to facilitate learning. As we point out
315 in this work, due to the opponent plasticity rules in dSPNs and iSPNs, a post-hoc efferent signal
316 that is not causally relevant to the decision-making process is necessary for appropriate learning.

317 Other authors have proposed models in which efferent activity is used for learning. In the context of
318 vocal learning in songbirds, Fee and Goldberg (2011) proposed that the variability-generating area
319 LMAN, which projects to the song motor pathway, sends collateral projections to Area X, which
320 undergoes dopamine-modulated plasticity. In this model, the efferent inputs to Area X allow it to
321 learn which motor commands are associated with better song performance (signaled by dopamine).
322 Similar to our model, this architecture implements off-policy reinforcement learning in Area X,
323 with HVC inputs to Area X being analogous to corticostriatal projections in our model. However,
324 in our work, the difference in plasticity rules between dSPNs and iSPNs is key to avoiding inter-
325 ference between efferent learning-related activity and feedforward action selection-related activity.
326 A similar architecture was proposed in Fee (2012) in the context of oculomotor learning, in which
327 oculomotor striatum receives efferent collaterals from the superior colliculus and/or cortical areas
328 which generate exploratory variability. Lisman (2014) also propose a high-level model of striatal
329 efferent inputs similar to ours, and also point out the issue with the iSPN plasticity rule assigning
330 credit to inappropriate actions without efferent inputs.

331 Our model is consistent with these prior proposals, but describes how efferent input must be
332 targeted to opponent SPNs. In our work, the distinction between dSPN and iSPN plasticity rules
333 is key to enable multiplexing of action-selection and efferent learning signals without interference.
334 Previous authors have proposed other mechanisms to avoid interference. For instance, Fee (2014)
335 propose that efferent inputs might influence plasticity without driving SPN spiking by synapsing
336 preferentially onto dendritic shafts rather than spines. To avoid action selection-related spikes
337 interfering with learning, the system may employ spike timing-dependent plasticity rules that are
338 tuned to match the latency at which efferent inputs excite SPNs. While these hypotheses are
339 not mutually exclusive to ours, our model requires no additional circuitry or assumptions beyond
340 the presence of appropriately tuned efferent input (see below) and opposite-sign plasticity rules
341 in dSPNs and iSPNs, due to the orthogonality of the sum and difference modes. An important
342 capability enabled by our model is that action selection and efferent inputs can be multiplexed
343 simultaneously, unlike the works cited above, which posit the existence of temporally segregated
344 action-selection and learning phases of SPN activity.

345 **Biological substrates of striatal efferent inputs**

346 Efferent inputs to the striatum must satisfy two important conditions for our model to learn cor-
347 rectly. First, they must be appropriately targeted: when an action is performed, dSPNs and
348 iSPNs associated with that action must be excited, but other dSPNs and iSPNs must not be. The
349 striatum receives topographically organized inputs from cortex (Peters et al., 2021) and thalamus
350 (Smith et al., 2004), and SPNs tuned to the same behavior tend to be located nearby in space
351 (Barbera et al., 2016; Shin et al., 2020; Klaus et al., 2017). This anatomical organization could
352 enable action-specific efferent inputs. Another possibility is that targeting of efferent inputs could
353 be tuned via plasticity during development. For instance, if a dSPN promotes a particular action,
354 reward-independent Hebbian plasticity of its efferent inputs would potentiate those inputs that en-
355 code the promoted action. Reward-independent anti-Hebbian plasticity would serve an analogous
356 function for iSPNs. Alternatively, if efferent inputs are fixed, plasticity downstream of striatum
357 could adapt the causal effect of SPNs to match their corresponding efferent input.

358 A second key requirement of our model is that efferent input synapses should not be adjusted
359 according to the same reward-modulated plasticity rules as the feedforward corticostriatal inputs,
360 as these rules would disrupt the targeting of efferent inputs to the corresponding SPNs. This
361 may be achieved in multiple ways. One possibility is that efferent inputs project from different
362 subregions or cell types than feedforward inputs and are subject to different forms of plasticity.
363 Alternatively, efferent input synapses may have been sufficiently reinforced that they exist in a less
364 labile, “consolidated” synaptic state. A third possibility is that the system may take advantage of
365 latency in efferent activity. Spike timing dependence in SPN input plasticity has been observed in
366 several studies (Shen et al., 2008; Fino et al., 2005; Pawlak and Kerr, 2008; Fisher et al., 2017).
367 This timing dependence could make plasticity sensitive to paired activity in state inputs and SPNs
368 while being insensitive to paired activity in efferent inputs and SPNs. Investigating the source of
369 efferent inputs to SPNs and how it is differentiated from other inputs is an important direction for
370 future work.

371 **Extensions and future work**

372 We have assumed that the striatum selects among a finite set of actions, each of which corresponds
373 to mutually uncorrelated patterns of SPN activity. In reality, there is evidence that the striatal
374 code for action is organized such that kinematically similar behaviors are encoded by similar SPN
375 activity patterns (Klaus et al., 2017; Markowitz et al., 2018). Other work has shown that the
376 dorsolateral striatum can exert influence over detailed kinematics of learned motor behaviors, rather
377 than simply select among categorically distinct actions (Dhawale et al., 2021). A more continuous,
378 structured code for action in dorsolateral striatum is useful in allowing reinforcement learning
379 to generalize between related actions. The ability afforded by our model to multiplex arbitrary
380 action selection and learning signals may facilitate these more sophisticated coding schemes. For
381 instance, reinforcement learning in continuous-valued action spaces requires a three-factor learning
382 rule in which the postsynaptic activity factor represents the discrepancy between the selected action
383 and the action typically selected in the current behavioral state (Lindsey and Litwin-Kumar, 2022),
384 which in our model would be represented by efferent activity in SPNs. Investigating such extensions
385 to our model and their consequences for SPN tuning is an interesting future direction.

386 In this work we find strong empirical evidence for our model of efferent activity in SPNs and
387 show that in principle it enables off-policy reinforcement learning capabilities. A convincing ex-
388 perimental demonstration of off-policy learning capabilities would require a way of identifying the
389 causal contribution of SPN activity to action selection, in order to distinguish between actions that
390 are consistent (on-policy) or inconsistent (off-policy) with SPN outputs. This could be achieved
391 through targeted stimulation of SPN populations, or by recording SPN activity during behaviors
392 that are known to be independent of striatal influence (Mizes et al., 2023). Simultaneous record-
393 ings in SPNs and other brain regions would also facilitate distinguishing between actions driven by
394 striatum from those driven by other pathways. Our model predicts that the relative strength of
395 fluctuations in difference mode versus sum mode activity should be greatest during striatum-driven
396 actions. Such experimental design would also enable a stronger test of the Q-learning model of
397 dopamine activity: actions driven by other regions should lead to increased dopamine activity, as
398 they will be predicted according to the striatum’s learned action-values to have low value.

399 In our model, the difference between dSPN and iSPN plasticity rules is key to enabling multiplexing
400 of action-selection and learning-related activity without interference. Observed plasticity rules
401 elsewhere in the brain are also heterogeneous; for instance, both Hebbian and anti-Hebbian behavior
402 are observed in cortico-cortical connections (Koch et al., 2013; Chindemi et al., 2022). It is an
403 interesting question whether a similar strategy may be employed outside the striatum, and in other
404 contexts besides reinforcement learning, to allow simultaneous encoding of behavior and learning-
405 related signals without interference.

406 Acknowledgments

407 We thank Jaeon Lee for providing the initial inspiration for this project, Sean Escola for fruitful
408 discussions, and Steven A. Siegelbaum for comments on the manuscript. J.L. is supported by the
409 Mathers Foundation and the Gatsby Charitable Foundation. J.M. is supported by a Career Award
410 at the Scientific Interface from the Burroughs Wellcome Fund, a fellowship from the Sloan Foun-
411 dation, and a fellowship from the David and Lucille Packard Foundation. S.R.D. is supported by
412 NIH grants RF1AG073625, R01NS114020, U24NS109520, the Simons Foundation Autism Research
413 Initiative, and the Simons Collaboration on Plasticity and the Aging Brain. A.L.-K. is supported
414 by the Mathers Foundation, the Burroughs Wellcome Foundation, the McKnight Endowment Fund,
415 and the Gatsby Charitable Foundation.

416 Declaration of interests

417 S.R.D. sits on the scientific advisory boards of Neumora and Gilgamesh Therapeutics, which have
418 licensed or sub-licensed the MoSeq technology.

419 Methods

420 Numerical simulations

421 Basic model architecture

422 In our simulated learning tasks, we used networks with the following architecture.

423 SPNs receive inputs from cortical neurons. In our simulated go/no-go tasks, there is a single cortical
424 input neuron (representing a task cue) with activity equal to 1 on each trial. In simulated tasks with
425 multiple different task cues (such as the two-alternative forced choice task), there is a population
426 of cortical input neurons, each of which is active with activity 1 when the corresponding task cue
427 is presented and 0 otherwise. The task cue is randomly chosen with uniform probability each trial.

428 For each of the A actions available to the model, there is an assigned dSPN and iSPN. We choose to
429 use a single neuron per action for simplicity of the model, but our model could easily be generalized
430 to use population activity to encode actions. The activities of the dSPN and iSPN associated with
431 action a are denoted as y_a^{dSPN} and y_a^{iSPN} , respectively. Each dSPN and iSPN receives inputs from M
432 cortical neurons, and the synaptic input weights from cortical neuron j to dSPN or iSPN associated
433 with action a are denoted as w_{aj}^{dSPN} or w_{aj}^{iSPN} . Feedforward SPN activity is given by

$$y_a^{\text{dSPN}} = \phi \left(\sum_{j=1}^M w_{aj}^{\text{dSPN}} x_j \right) \quad (8)$$

$$y_a^{\text{iSPN}} = \phi \left(\sum_{j=1}^M w_{aj}^{\text{iSPN}} x_j \right) \quad (9)$$

434 where ϕ is a nonlinear activation function. We choose ϕ to be the rectified linear function: $\phi(h) =$
435 $\max(0, h)$.

436 Action selection depends on SPN activity in the following manner. The log-likelihood of an action
437 a being performed is proportional to $\ell_a = y_a^{\text{dSPN}} - y_a^{\text{iSPN}}$. That is, dSPN activity increases the
438 likelihood of taking the action and iSPN activity decreases the likelihood of taking the action.
439 Concretely, the probability of action a being taken is:

$$p(a) = \frac{e^{\beta \ell_a}}{c_{\text{no-go}} + \sum_{a'} e^{\beta \ell_{a'}}} \quad (10)$$

440 where β is a parameter controlling the degree of stochasticity in action selection (higher β corre-
441 sponds to more deterministic choices), and c controls the probability that no action is taken. In
442 the simulated go/no-go tasks we choose $c_{\text{no-go}} = 1$ and in the tasks involving selection among
443 multiple actions we choose $c_{\text{no-go}} = 0$. Except where otherwise noted we used $\beta = 10.0$ in all task
444 simulations.

445 Models of SPN activity following action selection

446 In the “canonical action selection model” (Fig. 1), following action selection, the activity of the
447 dSPN associated with the selected action and the activity of all iSPNs associated with unselected
448 actions are set to 1. Biologically, this activity pattern can be implemented via effective mutual
449 inhibition between SPNs with opponent functions (dSPNs tuned to different actions, iSPNs tuned
450 to different actions, and dSPN/iSPN pairs tuned to the same action) and mutual excitation between
451 SPNs with complementary functions (dSPNs tuned to one action and iSPNs to another) (Burke
452 et al., 2017).

453 In the proposed efference model, following selection of an action a^* , activity of the SPNs associated
454 with action a^* is updated as follows:

$$y_a^{\text{dSPN}} \leftarrow \phi \left(c_{\text{efference}} \cdot 1[a = a^*] + \sum_{j=1}^M w_{aj}^{\text{dSPN}} x_j \right) \quad (11)$$

$$y_a^{\text{iSPN}} \leftarrow \phi \left(c_{\text{efference}} \cdot 1[a = a^*] + \sum_{j=1}^M w_{aj}^{\text{iSPN}} x_j \right) \quad (12)$$

$$(13)$$

455 where $1[a = a^*]$ equals 1 for $a = a^*$ and 0 otherwise. The parameter c controls the strength of
456 efferent excitation.

457 Learning rules

458 In all models, SPN input weights are initialized at 1 and weight updates proceed according to the
459 plasticity rules given below:

$$\Delta w_{aj}^{\text{dSPN}} = \alpha \left(f^{\text{dSPN}}(\delta) \cdot y_a^{\text{dSPN}} \cdot x_j \right), \quad (14)$$

$$\Delta w_{aj}^{\text{iSPN}} = \alpha \left(f^{\text{iSPN}}(\delta) \cdot y_a^{\text{iSPN}} \cdot x_j \right), \quad (15)$$

460 where α is a learning rate, set to 0.05 throughout all learning simulations (except the tutoring
461 simulations of Fig. 6 where it is set to 0.01). In the paper we experiment with various choices of
462 f^{dSPN} and f^{iSPN} .

$$f^{\text{dSPN}}(\delta) = \delta, f^{\text{iSPN}}(\delta) = -\delta \quad (\text{Linear}) \quad (16)$$

$$f^{\text{dSPN}}(\delta) = \max(\delta, 0), f^{\text{iSPN}}(\delta) = \max(-\delta, 0) \quad (\text{Rectified}) \quad (17)$$

$$f^{\text{dSPN}}(\delta) = \frac{1}{2} \left(a + \left(\frac{b}{(1 + ce^{1-d\delta})} \right) \right), f^{\text{iSPN}}(\delta) = \frac{1}{2} \left(a + \left(\frac{b}{(1 + ce^{1+d\delta})} \right) \right) \quad (\text{Offset sigmoid}) \quad (18)$$

463 with the offset sigmoid parameters chosen as $a = -3.5, b = 11.5, c = 0.9, d = 1$ (taken from Cruz
464 et al. (2022)). The quantity δ indicates an estimate of reward prediction error. In our experiments
465 in Fig 2 and Fig. 3 we use temporal difference learning to compute δ :

$$\delta = r - V(s) \quad (19)$$

$$\Delta V(s) = \alpha_V \delta \quad (20)$$

466 where α_V is a learning rate, set to 0.05 throughout all learning simulations (except the tutoring
467 simulations of Fig. 6 where it is set to 0.25) and s indicates the cortical input state (indicating
468 which cue is being presented). $V(s)$ is initialized at 0.

469 In our experiments in Fig. 6 we use Q-learning to enable off-policy learning, corresponding to the
470 following value for δ :

$$\delta = r - Q(s, a) \tag{21}$$

$$\tag{22}$$

471 where a indicates the action that was just taken in response to state s , and $Q(s, a)$ is taken to be
472 equal to the striatal output $\ell_a = y_a^{\text{dSPN}} - y_a^{\text{iSPN}}$ in response to the state s .

473 **Experimental prediction simulations**

474 For the model predictions of Fig. 4 and Fig. 5, we used the following parameters: $A = 50$, $\beta =$
475 100 , $c_{\text{efference}} = 1.5$ and set $c_{\text{no-go}}$ such that the no-action option was chosen 50% of the time.
476 Feedforward SPN activity was generated from a Gaussian process with kernel $k(t_1, t_2) = e^{-|t_1 - t_2|/10}$
477 (exponentially decaying autocorrelation with a time constant of 10 timesteps). Efference activity
478 also decayed exponentially with a time constant of 10 timesteps. Action selection occurred every 10
479 timesteps based on the SPN activity at the preceding timestep.

480 **Neural data analysis**

481 For our analysis of SPN data we used recordings previously described by Markowitz et al. (2018).
482 For our analysis of dopamine data we used the recordings described in Markowitz et al. (2023).

483 **Fiber photometry data**

484 Adeno-associated viruses (AAVs) expressing Cre-On jRCaMP1b and Cre-Off GCaMP6s were in-
485 jected into the dorsolateral striatum (DLS) of $n = 10$ *Drd1a-Cre* mice to measure bulk dSPN (red)
486 and iSPN (green) activity via multicolor photometry. Activity of each indicator was recorded at
487 a rate of 30Hz using an optical fiber implanted in the right DLS. Data was collected during spon-
488 taneous behavior in a circular open field, for 5-6 sessions of 20 minutes each for each mouse. In
489 the reversed indicator experiments of Supp. Fig. 2, *A2a-Cre* mice were injected with a mixture of
490 the same AAVs, labeling iSPNs with jRCaMP1b (red) and dSPNs with GCaMP6s (green). More
491 details are reported in Markowitz et al. (2018).

492 In our data analyses in Fig. 4C and Supp. Fig 2, for each session ($n = 48$ and $n = 8$, respectively)
493 we computed the autocorrelation and cross-correlation of the dSPN and iSPN indicator activity
494 across the entire session.

495 **Miniscope data**

496 *Drd1a-Cre* AAVs expressing GCaMP6f were injected into the right DLS of $n = 4$ *Drd1a-Cre* mice (to
497 label dSPNs) and $n = 6$ *A2a-Cre* mice (to label iSPNs). A head-mounted single-photon microscope
498 was coupled to a gradient index lens implanted into the dorsal striatum above the injection site.
499 Recordings were made, as for the photometry data, during spontaneous behavior in a circular open
500 field. Calcium activity was recorded from a total of 653 dSPNs and 794 iSPNs for these mice, with
501 the number of neurons per mouse ranging from 27–336. To enable simultaneous recording of dSPNs
502 and iSPNs in the same mice, a different protocol was used: *Drd1a-Cre* mice were injected with an
503 AAV mixture which labeled both dSPNs and iSPNs with GCaMP6s, but additionally selectively
504 labeled dSPNs with nuclear-localized dTomato. This procedure enabled (in $n = 4$ mice) cell-type
505 identification of dSPNs vs. iSPNs with a two-photon microscope which was cross-referenced with
506 the single-photon microscope recordings. More details are given in Markowitz et al. (2018). In our
507 analyses, these data were used for the simultaneous-recording analyses in Fig. 5L,M,N,O and were
508 also combined with the appropriate single-pathway data in the analyses of Fig. 5J,K.

509 **Behavioral data**

510 Mouse behavior in the circular open field was recorded as follows: 3D pose information was recorded
511 using a depth camera at a rate of 30Hz. The videos were preprocessed to center the mouse and align
512 the nose-to-tail axis across frames and remove occluding objects. The videos were then fed through
513 PCA to reduce the dimensionality of the data and fed into the MoSeq algorithm (Wiltschko et al.,
514 2015) which fits a generative model to the video data that automatically infers a set of behavioral
515 “syllables” (repeated, stereotyped behavioral kinematics) and assigns each frame of the video to
516 one of these syllables. More details on MoSeq are given in Wiltschko et al. (2015) and more details
517 on its application to this dataset are given in Markowitz et al. (2018). There were 89 syllables
518 identified by MoSeq that appear across all the sessions. We restricted our analysis to the set of 62
519 syllables that appear at least 5 times in each behavioral session.

520 **Syllable-tuned SPN activity mode analysis**

521 In our analysis, we first z-scored the activity of each neuron across the data collected for each mouse.
522 We divided the data by the boundaries of behavioral syllables and split it into two equally sized
523 halves (based on whether the timestamp, rounded to the nearest second, of the behavioral syllable
524 was even or odd). To compute the activity modes associated with each behavioral syllable, we
525 first computed the average change in activity for each neuron during each syllable and fit a linear
526 regression model to predict this increase from a one-hot vector indicating the syllable identity.
527 The resulting coefficients of this regression indicate the directions (“modes”) in activity space that
528 increase the most during performance of each of the behavioral syllables. We linearly time-warped
529 the data in each session based on the boundaries of each MoSeq-identified behavioral syllable, such
530 that in the new time coordinates each behavioral syllable lasted 10 timesteps. The time course of
531 the projection of SPN activity along the modes associated with each behavioral syllable was then
532 computed around the onset of that syllable, or around all other syllables. As a way of crossvalidating
533 the analysis, we performed the regression on one half of the data and plotted the average mode
534 activity on the other half of the data (in both directions, and averaged the results). We averaged

535 the resulting time courses of mode activity across all choices of behavioral syllables. This analysis
536 was performed for each mouse and the results in Fig. 5J,K,L,M,N,O show means and standard
537 errors across mice.

538 Dopamine activity data and analysis

539 For 6E we used data from Markowitz et al. (2023). Mice ($n = 14$) virally expressing the dopamine
540 reporter dLight1.1 in the DLS were recorded with a fiber cannula implanted above the injection
541 site. Mice were placed in a circular open field for 30 minute sessions and allowed to behave freely
542 while spontaneous dLight activity was recorded. MoSeq (described above) was used to infer a set
543 of $S57$ behavioral syllables observed across all sessions. As in Markowitz et al. (2023), the data
544 were preprocessed by computing the maximum dLight value during each behavioral syllable. These
545 per-syllable dopamine values were z-scored across each session and used as our measure of dopamine
546 activity during each syllable. We then computed an $S \times S$ table of the average dopamine activity
547 during each syllable s_t conditioned on the previous syllable having been syllable s_{t-1} , denoted as
548 $D(s_{t-1}, s_t)$. We also computed the $S \times X$ table of probabilities of transitioning from syllable s' to
549 syllable s across the dataset, denoted as $P(s_{t-1}, s_t)$. These tables were computed separately for
550 each mouse. In Fig. 6E we report the pearson correlation coefficient between the predicted and
551 actual values of $P(s_{t-1}, s_t)$. We then experimented with several alternative models (see Supple-
552 mental Information) that predict $P(s_{t-1}, s_t)$ based on $D(s_{t-1}, s_t)$. In Fig. 6E we report the pearson
553 correlation coefficient between the predicted and actual values of $P(s_{t-1}, s_t)$.

554

555 Supplemental information

556 Relationship between sum mode activity and future difference mode activity

557 In the main text we provided an argument for why sum mode activity drives changes to future
 558 difference mode activity, assuming a linear $f^{d/iSPN}(\delta)$ and linear neural activation functions. Here
 559 we generalize this argument to more general learning rules and activation functions ϕ , assuming
 560 only that $f^{dSPN}(\delta)$ is monotonically increasing, $f^{iSPN}(\delta)$ is monotonically increasing, and $\phi(\cdot)$ is
 561 monotonically increasing. We have that $y^{d/iSPN} = \phi(\mathbf{w}^{d/iSPN} \cdot \mathbf{x})$, and $\delta \mathbf{w}^{d/iSPN} = (f^{d/iSPN}(\delta) \cdot$
 562 $y^{d/iSPN})\mathbf{x}$. Thus, in the limit of small small weight updates, we can write:

$$\begin{aligned}
 \Delta(y^{dSPN} - y^{iSPN}) &= \Delta\phi(\mathbf{w}^{dSPN} \cdot \mathbf{x}) - \Delta\phi(\mathbf{w}^{iSPN} \cdot \mathbf{x}) \\
 &\approx \phi'(\mathbf{w}^{dSPN} \cdot \mathbf{x})(\Delta\mathbf{w}^{dSPN} \cdot \mathbf{x}) - \phi'(\mathbf{w}^{iSPN} \cdot \mathbf{x})(\Delta\mathbf{w}^{iSPN} \cdot \mathbf{x}) \\
 &\propto \phi'(\mathbf{w}^{dSPN} \cdot \mathbf{x})(f^{dSPN}(\delta) \cdot y^{dSPN} \mathbf{x} \cdot \mathbf{x}) - \phi'(\mathbf{w}^{iSPN} \cdot \mathbf{x})(f^{iSPN}(\delta) \cdot y^{iSPN} \mathbf{x} \cdot \mathbf{x}) \\
 &= \|\mathbf{x}\|^2 \left(\phi'(\mathbf{w}^{dSPN} \cdot \mathbf{x})(f^{dSPN}(\delta) \cdot y^{dSPN}) - \phi'(\mathbf{w}^{iSPN} \cdot \mathbf{x})(f^{iSPN}(\delta) \cdot y^{iSPN}) \right) \\
 &\propto c^{dSPN} f^{dSPN}(\delta) y^{dSPN} + (-c^{iSPN} f^{iSPN}(\delta) y^{iSPN}). \tag{23}
 \end{aligned}$$

563 where c^{dSPN} and c^{iSPN} are nonnegative because ϕ' is always nonnegative by assumption. Since by
 564 assumption $f^{d/iSPN}$ are increasing/decreasing, respectively, the first term of the above sum has
 565 nonnegative correlation with δy^{dSPN} and the second term has nonnegative correlation with δy^{iSPN} .
 566 Thus, changes $\Delta(y^{dSPN} - y^{iSPN})$ to difference mode activity are always nonnegatively correlated
 567 with sum mode activity. If we assume that efferent excitation is always sufficiently strong that
 568 $c^{dSPN} = \phi'(\mathbf{w}^{dSPN} \cdot \mathbf{x})$ and $c^{iSPN} = \phi'(\mathbf{w}^{iSPN} \cdot \mathbf{x})$ are positive, and that there are no values of δ
 569 for which $f^{d/iSPN}(\delta)$ both have zero derivative, we can further guarantee that changes to difference
 570 mode activity will always be *positively* correlated with sum mode activity.

571 Generalizing the model to a distributed code for actions

572 In our model simulations in the main text we assumed for convenience that there is a single dSPN
 573 and iSPN that promote and suppress each available action, respectively. It is more realistic to model
 574 the code for action as distributed among many SPNs. Our model generalizes easily to this case; all
 575 that is necessary is for the efferent activity following action selection to excite the vectors (for both
 576 dSPNs and iSPNs) in population activity space corresponding to that action. To demonstrate this,
 577 we conducted a simulation with $N = 1000$ dSPNs and iSPNs each, $S = 10$ input cues (one-hot
 578 input vectors), and $A = 10$ actions, with one correct action for each input state. Feedforward SPN
 579 activity is given by

$$y_i^{\text{dSPN}} = \phi \left(\sum_{j=1}^M w_{ij}^{\text{dSPN}} x_j \right) \quad (24)$$

$$y_i^{\text{iSPN}} = \phi \left(\sum_{j=1}^M w_{ij}^{\text{iSPN}} x_j \right) \quad (25)$$

580 The log-likelihood of an action a being performed is proportional to

$$\ell_a = \sum_{i=1}^N \zeta_{ai}^{\text{dSPN}} y_i^{\text{dSPN}} - \zeta_{ai}^{\text{iSPN}} y_i^{\text{iSPN}} \quad (26)$$

581 where ζ_{ai}^{dSPN} and ζ_{ai}^{iSPN} are randomly sampled uniformly in the interval $[0, 1]$ and then normalized
 582 so that each vector $\zeta_{\mathbf{a}}^{\text{dSPN}}$ and $\zeta_{\mathbf{a}}^{\text{iSPN}}$ has norm 1. Thus, the contribution of each dSPN/iSPN to
 583 the promotion/suppression of each action is randomly distributed.

584 In the efference model, following selection of an action a^* , activity of the SPNs associated with action
 585 a^* is updated as follows, so that efference activity excites the modes $\zeta_{\mathbf{a}^*}^{\text{dSPN}}$ and $\zeta_{\mathbf{a}^*}^{\text{iSPN}}$ associated
 586 with the selected action:

$$y_i^{\text{dSPN}} \leftarrow \phi \left(c_{\text{efference}} \cdot \zeta_{a^*i}^{\text{dSPN}} + \sum_{j=1}^M w_{ij}^{\text{dSPN}} x_j \right) \quad (27)$$

$$y_i^{\text{iSPN}} \leftarrow \phi \left(c_{\text{efference}} \cdot \zeta_{a^*i}^{\text{iSPN}} + \sum_{j=1}^M w_{ij}^{\text{iSPN}} x_j \right) \quad (28)$$

$$(29)$$

587 We also experiment with a generalization of the canonical action selection model to this distributed
 588 action tuning architecture, in which following action selection, SPN activity is set to

$$y_i^{\text{dSPN}} \leftarrow \zeta_{a^*i}^{\text{dSPN}} \quad (30)$$

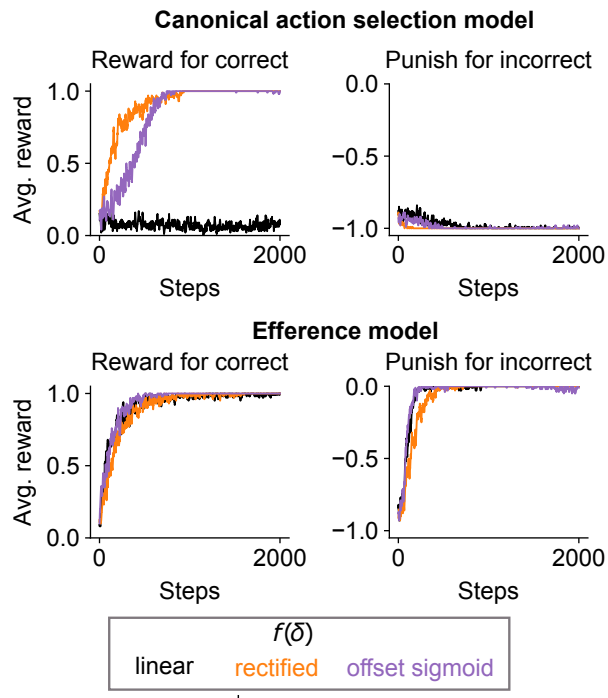
$$y_i^{\text{iSPN}} \leftarrow \left(\max_{i'} \zeta_{a^*i'}^{\text{iSPN}} \right) - \zeta_{a^*i}^{\text{iSPN}} \quad (31)$$

$$(32)$$

589 In this model, dSPNs are excited in proportion to their contribution to the currently selected action
 590 and iSPNs are suppressed in proportion to their degree of inhibition of the currently selected action.

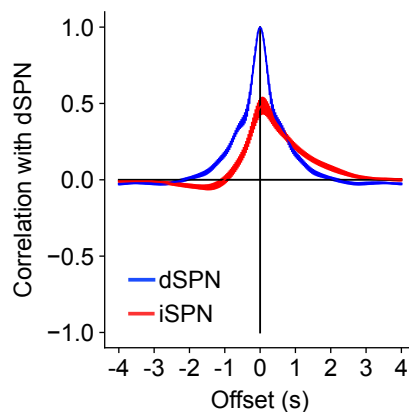
591 The plasticity rules used are the same as in the main text.

592 We find that the results of the main text – that the canonical action selection model fails to learn
593 from negative rewards, while the efference model successfully learns from both reward protocols –
594 is replicated (Supp. Fig. 1).



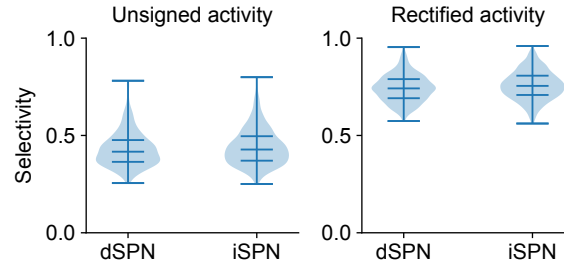
Supplemental Fig. 1: Performance of striatal RL models with a distributed code for actions on a task with 10 cortical input states, 10 available actions, and one correct action for each input state.

595 **Photometry analysis with reversed indicators**



Supplemental Fig. 2: Same as Fig. 4C, but performing the analysis on subjects with reversed assignment of indicators to SPN types.

596 **Comparison of selectivity of dSPNs and iSPNs**



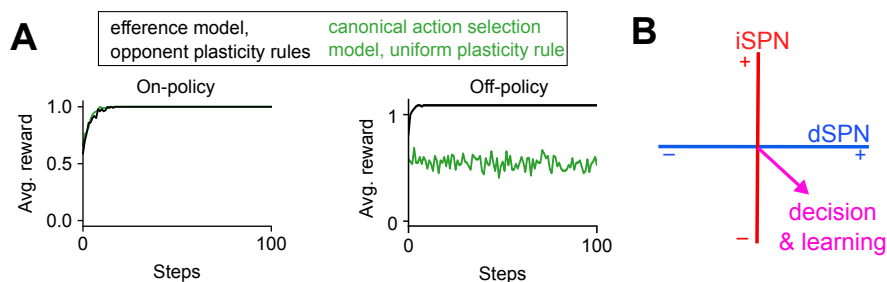
Supplemental Fig. 3: Comparison of dSPN and iSPN tuning selectivity. Violin plots indicate the distribution of selectivity values across all neurons computed using Eq. 33, using either unsigned (left) or rectified (right) z-scored activity as the raw measure of a neuron’s tuning to a behavioral syllable. Horizontal lines indicate the 0, 25, 50, 75, 100 percentile values of the distribution.

597 To test whether dSPNs or iSPNs exhibit greater or less specificity in their tuning to behaviors,
598 we computed the selectivity of each neuron in the imaging data of Fig. 5. For each neuron, we
599 computed its average z-scored activity a_i in response to each of the behavioral syllables $i \in \{1, \dots, A\}$
600 in the dataset. Common measures of selectivity require a nonnegative measurement of a neuron’s
601 tuning to a given condition. Thus, we conducted the analysis in two ways, using either the unsigned
602 activity $|a_i|$ or the rectified activity $\max(a_i, 0)$ as the measure of the neuron’s tuning t_i to syllable i .
603 The selectivity was then computed using the following expression introduced in prior work (Treves
604 and Rolls, 1991; Willmore and Tolhurst, 2001):

$$\frac{\left(\frac{1}{A} \sum_i t_i\right)^2}{\frac{1}{A} \sum_i t_i^2} \quad (33)$$

605 This value ranges from 0 to 1, and higher value indicates that fluctuations in a neuron’s activity are
606 driven primarily by one or a few behavioral syllables. The results are shown in Supp. Fig. 3. The
607 selectivity values are fairly modest (consistent with a distributed code for actions) and comparable
608 between dSPNs and iSPNs.

609 Alternative model with shared plasticity rule among all SPNs



Supplemental Fig. 4: Comparison to counterfactual model in which iSPNs use the same plasticity rule as dSPNs. A. Left: performance of simulated striatal RL system using efference model with the opponent dSPN/iSPN plasticity rules used elsewhere in the paper (black, same as Fig. 3E), and a system using the canonical action selection model and identical dSPN and iSPN plasticity rules (green). Right: same as left panel, but in an off-policy setting in which another pathway controls behavior during and always chooses the correct action, and the performance of the striatal RL system is evaluated over time. Here the Q-learning model of dopamine activity is used. B. In the counterfactual model in which iSPNs use the same plasticity rule as dSPNs, activity in the difference mode (dSPN - iSPN) influences (via plasticity) changes in future difference mode activity that affect decision-making.

610 The issues identified in Fig. 2 with the canonical action selection model are a consequence of the
611 iSPN plasticity rule. From a normative perspective is interesting to consider why the empirically
612 observed iSPN plasticity rule might be advantageous, compared to an alternative model in which
613 iSPNs share the same plasticity rule as dSPNs. For instance, this alternative model can solve
614 the two-alternative forced choice task of Fig. 2 with both positive and negative reward protocols
615 (Supp. Fig. 4A, left). However, the limitations of this alternative model are revealed in the off-
616 policy learning setting, where the Q-learning algorithm is required. In this case, SPN activity must
617 encode Q-values associated with each action, but in the canonical action selection model, these
618 values are disrupted by the updates to SPN activity following action selection. This is because
619 the activity updates in the canonical action selection model modify difference mode activity, which
620 (when dSPN and iSPN plasticity rules are the same) is needed for learning (Supp. Fig. 4B). As a
621 result, the predicted Q-values are inaccurate, and the model has difficulty learning the true value
622 of each action. We demonstrate this in the two-alternative forced task in an off-policy learning
623 protocol where an oracle chooses the correct action on each trial, and the striatal pathway's ability
624 to solve the task independently is evaluated. The efference activity model has no issue due to the
625 orthogonality of the efferent activity and difference modes as described above, but the canonical
626 action selection model fails to solve the task (Supp. Fig. 4A, right).

627 We note that non-orthogonality of the activity mode used for learning and behavior could cause
628 other problems besides impairing the system's ability to implement off-policy learning algorithms;
629 for instance, even in an on-policy setting it could interfere with sequential action selection at rapid
630 timescales.

631 Models used for dopamine analysis

632 We experimented with models that predict transition probabilities $P(s_{t-1}, s_t)$ based on average
633 dopamine activity $D(s_{t-1}, s_t)$ associated with each transition.

634

635 *Q-learning model:* In the Q-learning model, the mouse maintains an internal estimate of the value
636 $Q(s_{t-1}, s_t)$ of each transition between syllables. In the absence of explicit rewards, the dopamine
637 activity associated with a syllable transition is predicted to be: $D(s_{t-1}, s_t) = \max_{s'} Q(s_t, s') -$
638 $Q(s_{t-1}, s_t)$. We inferred a set of Q-values by initializing a Q-table with all zero values and running
639 gradient descent on the Q-table to minimize the mean squared error between the predicted and
640 empirical values of $D(s_{t-1}, s_t)$. These inferred Q-values were used to predict behavioral transition
641 probabilities according to: $\hat{P}(s_{t-1}, s_t) = \frac{e^{\beta(s_{t-1})Q(s_{t-1}, s_t)}}{\sum_{s'} e^{\beta(s_{t-1})Q(s_{t-1}, s')}}$. We did not fit the value of $\beta(s_{t-1})$ but
642 rather chose it to be the reciprocal of the standard deviation of $Q(s_{t-1}, s')$ across all s' , to ensure
643 a reasonable dynamic range in predicted transition probabilities.

644 *TD learning model:* In this model, the mouse maintains an internal estimate of the value $V(s)$
645 of each syllable, and the predicted dopamine activity at each transition is $D(s_{t-1}, s_t) = V(s_t) -$
646 $V(s_{t-1})$. We fit the vector of values $V(s)$ to minimize the mean squared error of predicted and
647 empirical $D(s_{t-1}, s_t)$. The predicted transition probabilities in this model (which are independent
648 of the previous syllable s_{t-1}) are: $\hat{P}(s_{t-1}, s_t) = \frac{e^{\beta V(s_t)}}{\sum_{s'} e^{\beta V(s')}}$ with β chosen to normalize the $V(s')$ to
649 have standard deviation 1, as in the previous models.

650 *Action value model:* In this model, we assume that dopamine activity simply reflects the proba-
651 bility of each transition rather than encoding a prediction error; that is, we assume $P(s_{t-1}, s_t) =$
652 $\frac{D(s_{t-1}, s_t)}{\sum_s D(s_{t-1}, s)}$.

653 *State value model:* In this model, we assume that dopamine activity simply reflects the proba-
654 bility of each behavioral syllable being chosen and is independent of the previous syllable. That
655 is, we compute the average dopamine activity $D(s)$ associated with each syllable s , and predict
656 $P(s_{t-1}, s_t) = \frac{D(s_t)}{\sum_s D(s)}$.

657 References

- 658 Arulkumaran, K., Deisenroth, M. P., Brundage, M., and Bharath, A. A. (2017). Deep reinforcement
659 learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6):26–38.
- 660 Ashby, F. G., Turner, B. O., and Horvitz, J. C. (2010). Cortical and basal ganglia contributions to
661 habit learning and automaticity. *Trends in cognitive sciences*, 14(5):208–215.
- 662 Balleine, B. W., Delgado, M. R., and Hikosaka, O. (2007). The role of the dorsal striatum in reward
663 and decision-making. *Journal of Neuroscience*, 27(31):8161–8165.
- 664 Barbera, G., Liang, B., Zhang, L., Gerfen, C. R., Culurciello, E., Chen, R., Li, Y., and Lin, D.-
665 T. (2016). Spatially compact neural clusters in the dorsal striatum encode locomotion relevant
666 information. *Neuron*, 92(1):202–213.
- 667 Bostan, A. C. and Strick, P. L. (2018). The basal ganglia and the cerebellum: nodes in an integrated
668 network. *Nature Reviews Neuroscience*, 19(6):338–350.
- 669 Burke, D. A., Rotstein, H. G., and Alvarez, V. A. (2017). Striatal local circuitry: a new framework
670 for lateral inhibition. *Neuron*, 96(2):267–284.
- 671 Calabresi, P., Gubellini, P., Centonze, D., Picconi, B., Bernardi, G., Chergui, K., Svenningsson,
672 P., Fienberg, A. A., and Greengard, P. (2000). Dopamine and camp-regulated phosphoprotein
673 32 kda controls both striatal long-term depression and long-term potentiation, opposing forms
674 of synaptic plasticity. *Journal of Neuroscience*, 20(22):8443–8451.
- 675 Cardinal, R. N., Parkinson, J. A., Hall, J., and Everitt, B. J. (2002). Emotion and motivation:
676 the role of the amygdala, ventral striatum, and prefrontal cortex. *Neuroscience & Biobehavioral*
677 *Reviews*, 26(3):321–352.
- 678 Chindemi, G., Abdellah, M., Amsalem, O., Benavides-Piccione, R., Delattre, V., Doron, M., Ecker,
679 A., Jaquier, A. T., King, J., Kumbhar, P., et al. (2022). A calcium-based plasticity model for
680 predicting long-term potentiation and depression in the neocortex. *Nature Communications*,
681 13(1):3038.
- 682 Collins, A. G. and Frank, M. J. (2014). Opponent actor learning (opal): modeling interactive
683 effects of striatal dopamine on reinforcement learning and choice incentive. *Psychological review*,
684 121(3):337.
- 685 Contreras-Vidal, J. L. and Schultz, W. (1999). A predictive reinforcement model of dopamine
686 neurons for learning approach behavior. *Journal of computational neuroscience*, 6(3):191–214.
- 687 Cruz, B. F., Guiomar, G., Soares, S., Motiwala, A., Machens, C. K., and Paton, J. J. (2022). Action
688 suppression reveals opponent parallel control via striatal circuits. *Nature*, 607(7919):521–526.
- 689 Cui, G., Jun, S. B., Jin, X., Pham, M. D., Vogel, S. S., Lovinger, D. M., and Costa, R. M. (2013).
690 Concurrent activation of striatal direct and indirect pathways during action initiation. *Nature*,
691 494(7436):238–242.
- 692 Dhawale, A. K., Wolff, S. B., Ko, R., and Ölveczky, B. P. (2021). The basal ganglia control the
693 detailed kinematics of learned motor skills. *Nature neuroscience*, 24(9):1256–1269.

- 694 Dreyer, J. K., Herrik, K. F., Berg, R. W., and Hounsgaard, J. D. (2010). Influence of phasic and
695 tonic dopamine release on receptor activation. *Journal of Neuroscience*, 30(42):14273–14283.
- 696 Exner, C., Koschack, J., and Irle, E. (2002). The differential role of premotor frontal cortex and
697 basal ganglia in motor sequence learning: evidence from focal basal ganglia lesions. *Learning &
698 Memory*, 9(6):376–386.
- 699 Fee, M. S. (2012). Oculomotor learning revisited: a model of reinforcement learning in the basal
700 ganglia incorporating an efference copy of motor actions. *Frontiers in neural circuits*, 6:38.
- 701 Fee, M. S. (2014). The role of efference copy in striatal learning. *Current opinion in neurobiology*,
702 25:194–200.
- 703 Fee, M. S. and Goldberg, J. H. (2011). A hypothesis for basal ganglia-dependent reinforcement
704 learning in the songbird. *Neuroscience*, 198:152–170.
- 705 Fino, E., Glowinski, J., and Venance, L. (2005). Bidirectional activity-dependent plasticity at
706 corticostriatal synapses. *Journal of Neuroscience*, 25(49):11279–11287.
- 707 Fisher, S. D., Robertson, P. B., Black, M. J., Redgrave, P., Sagar, M. A., Abraham, W. C.,
708 and Reynolds, J. N. (2017). Reinforcement determines the timing dependence of corticostriatal
709 synaptic plasticity in vivo. *Nature communications*, 8(1):334.
- 710 Frank, M. J. (2005). Dynamic dopamine modulation in the basal ganglia: a neurocomputational
711 account of cognitive deficits in medicated and nonmedicated parkinsonism. *Journal of cognitive
712 neuroscience*, 17(1):51–72.
- 713 Freeze, B. S., Kravitz, A. V., Hammack, N., Berke, J. D., and Kreitzer, A. C. (2013). Control of
714 basal ganglia output by direct and indirect pathway projection neurons. *Journal of Neuroscience*,
715 33(47):18531–18539.
- 716 Gurney, K. N., Humphries, M. D., and Redgrave, P. (2015). A new framework for cortico-striatal
717 plasticity: behavioural theory meets in vitro data at the reinforcement-action interface. *PLoS
718 biology*, 13(1):e1002034.
- 719 Houk, J. C. and Adams, J. L. (1995). 13 a model of how the basal ganglia generate and use neural
720 signals that. *Models of information processing in the basal ganglia*, page 249.
- 721 Hwang, E. J., Dahlen, J. E., Hu, Y. Y., Aguilar, K., Yu, B., Mukundan, M., Mitani, A., and
722 Komiyama, T. (2019). Disengagement of motor cortex from movement control during long-term
723 learning. *Science advances*, 5(10):eaay0001.
- 724 Iino, Y., Sawada, T., Yamaguchi, K., Tajiri, M., Ishii, S., Kasai, H., and Yagishita, S. (2020).
725 Dopamine d2 receptors in discrimination learning and spine enlargement. *Nature*, 579(7800):555–
726 560.
- 727 Ito, M. and Doya, K. (2011). Multiple representations and algorithms for reinforcement learning
728 in the cortico-basal ganglia circuit. *Current opinion in neurobiology*, 21(3):368–373.
- 729 Jaskir, A. and Frank, M. J. (2023). On the normative advantages of dopamine and striatal oppo-
730 nency for learning and choice. *Elife*, 12:e85107.
- 731 Joel, D., Niv, Y., and Ruppin, E. (2002). Actor–critic models of the basal ganglia: New anatomical
732 and computational perspectives. *Neural networks*, 15(4-6):535–547.

- 733 Kawai, R., Markman, T., Poddar, R., Ko, R., Fantana, A. L., Dhawale, A. K., Kampff, A. R., and
734 Ölviczky, B. P. (2015). Motor cortex is required for learning but not for executing a motor skill.
735 *Neuron*, 86(3):800–812.
- 736 Klaus, A., Martins, G. J., Paixao, V. B., Zhou, P., Paninski, L., and Costa, R. M. (2017). The
737 spatiotemporal organization of the striatum encodes action space. *Neuron*, 95(5):1171–1180.
- 738 Koch, G., Ponzio, V., Di Lorenzo, F., Caltagirone, C., and Veniero, D. (2013). Hebbian and
739 anti-hebbian spike-timing-dependent plasticity of human cortico-cortical connections. *Journal of*
740 *Neuroscience*, 33(23):9725–9733.
- 741 Kravitz, A. V., Freeze, B. S., Parker, P. R., Kay, K., Thwin, M. T., Deisseroth, K., and Kreitzer,
742 A. C. (2010). Regulation of parkinsonian motor behaviours by optogenetic control of basal ganglia
743 circuitry. *Nature*, 466(7306):622–626.
- 744 Lee, J. and Sabatini, B. L. (2021). Striatal indirect pathway mediates exploration via collicular
745 competition. *Nature*, 599(7886):645–649.
- 746 Lee, S. J., Lodder, B., Chen, Y., Patriarchi, T., Tian, L., and Sabatini, B. L. (2021). Cell-type-
747 specific asynchronous modulation of pka by dopamine in learning. *Nature*, 590(7846):451–456.
- 748 Lindsey, J. and Litwin-Kumar, A. (2022). Action-modulated midbrain dopamine activity arises
749 from distributed control policies. *Advances in Neural Information Processing Systems*, 35:5535–
750 5548.
- 751 Lisman, J. (2014). Two-phase model of the basal ganglia: implications for discontinuous control
752 of the motor system. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
753 369(1655):20130489.
- 754 Markowitz, J. E., Gillis, W. F., Beron, C. C., Neufeld, S. Q., Robertson, K., Bhagat, N. D.,
755 Peterson, R. E., Peterson, E., Hyun, M., Linderman, S. W., et al. (2018). The striatum organizes
756 3d behavior via moment-to-moment action selection. *Cell*, 174(1):44–58.
- 757 Markowitz, J. E., Gillis, W. F., Jay, M., Wood, J., Harris, R. W., Cieszkowski, R., Scott, R., Brann,
758 D., Koveal, D., Kula, T., et al. (2023). Spontaneous behaviour is structured by reinforcement
759 without explicit reward. *Nature*, 614(7946):108–117.
- 760 Mink, J. W. (1996). The basal ganglia: focused selection and inhibition of competing motor
761 programs. *Progress in neurobiology*, 50(4):381–425.
- 762 Mizes, K. G., Lindsey, J., Escola, G. S., and Ölviczky, B. P. (2023). Dissociating the contributions
763 of sensorimotor striatum to automatic and visually guided motor sequences. *Nature Neuroscience*,
764 pages 1–14.
- 765 Montague, P. R., Dayan, P., and Sejnowski, T. J. (1996). A framework for mesencephalic dopamine
766 systems based on predictive hebbian learning. *Journal of neuroscience*, 16(5):1936–1947.
- 767 Niv, Y. (2009). Reinforcement learning in the brain. *Journal of Mathematical Psychology*,
768 53(3):139–154.
- 769 O’Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., and Dolan, R. J. (2004). Dissocia-
770 ble roles of ventral and dorsal striatum in instrumental conditioning. *science*, 304(5669):452–454.

- 771 Packard, M. G. and Knowlton, B. J. (2002). Learning and memory functions of the basal ganglia.
772 *Annual review of neuroscience*, 25(1):563–593.
- 773 Pawlak, V. and Kerr, J. N. (2008). Dopamine receptor activation is required for corticostriatal
774 spike-timing-dependent plasticity. *Journal of Neuroscience*, 28(10):2435–2446.
- 775 Peak, J., Chieng, B., Hart, G., and Balleine, B. W. (2020). Striatal direct and indirect pathway
776 neurons differentially control the encoding and updating of goal-directed learning. *Elife*, 9:e58544.
- 777 Peters, A. J., Fabre, J. M., Steinmetz, N. A., Harris, K. D., and Carandini, M. (2021). Striatal
778 activity topographically reflects cortical activity. *Nature*, 591(7850):420–425.
- 779 Redgrave, P., Prescott, T. J., and Gurney, K. (1999). The basal ganglia: a vertebrate solution to
780 the selection problem? *Neuroscience*, 89(4):1009–1023.
- 781 Schultz, W., Dayan, P., and Montague, P. R. (1997). A neural substrate of prediction and reward.
782 *Science*, 275(5306):1593–1599.
- 783 Seo, M., Lee, E., and Averbeck, B. B. (2012). Action selection and action value in frontal-striatal
784 circuits. *Neuron*, 74(5):947–960.
- 785 Shen, W., Flajolet, M., Greengard, P., and Surmeier, D. J. (2008). Dichotomous dopaminergic
786 control of striatal synaptic plasticity. *Science*, 321(5890):848–851.
- 787 Shin, J. H., Song, M., Paik, S.-B., and Jung, M. W. (2020). Spatial organization of functional clus-
788 ters representing reward and movement information in the striatal direct and indirect pathways.
789 *Proceedings of the National Academy of Sciences*, 117(43):27004–27015.
- 790 Silveri, M. C. (2021). Contribution of the cerebellum and the basal ganglia to language production:
791 Speech, word fluency, and sentence construction—evidence from pathology. *The Cerebellum*,
792 20(2):282–294.
- 793 Smith, Y., Raju, D. V., Pare, J.-F., and Sidibe, M. (2004). The thalamostriatal system: a highly
794 specific network of the basal ganglia circuitry. *Trends in neurosciences*, 27(9):520–527.
- 795 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- 796 Treves, A. and Rolls, E. T. (1991). What determines the capacity of autoassociative memories in
797 the brain? *Network: Computation in Neural Systems*, 2(4):371.
- 798 Varin, C., Cornil, A., Houtteman, D., Bonnavion, P., and de Kerchove d’Exaerde, A. (2023).
799 The respective activation and silencing of striatal direct and indirect pathway neurons support
800 behavior encoding. *Nature communications*, 14(1):4982.
- 801 Wickens, J., Begg, A., and Arbuthnott, G. (1996). Dopamine reverses the depression of rat corti-
802 costriatal synapses which normally follows high-frequency stimulation of cortex in vitro. *Neuro-
803 science*, 70(1):1–5.
- 804 Wildgruber, D., Ackermann, H., and Grodd, W. (2001). Differential contributions of motor cortex,
805 basal ganglia, and cerebellum to speech motor control: effects of syllable repetition rate evaluated
806 by fmri. *Neuroimage*, 13(1):101–109.
- 807 Willmore, B. and Tolhurst, D. J. (2001). Characterizing the sparseness of neural codes. *Network:
808 Computation in Neural Systems*, 12(3):255.

809 Wiltschko, A. B., Johnson, M. J., Iurilli, G., Peterson, R. E., Katon, J. M., Pashkovski, S. L.,
810 Abraira, V. E., Adams, R. P., and Datta, S. R. (2015). Mapping sub-second structure in mouse
811 behavior. *Neuron*, 88(6):1121–1135.