

RESEARCH ARTICLE

# How cortico-basal ganglia-thalamic subnetworks can shift decision policies to increase reward rate

Jyotika Bahuguna<sup>1‡\*</sup>, Timothy Verstynen<sup>1,2‡\*</sup>, Jonathan E. Rubin<sup>2,3‡\*</sup>

**1** Department of Psychology and Neuroscience Institute, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States of America, **2** Center for the Neural Basis of Cognition, Pittsburgh, Pennsylvania, United States of America, **3** Department of Mathematics, University of Pittsburgh, Pittsburgh, Pennsylvania, United States of America

‡ Current address: Université de Strasbourg, Laboratoire de Neurosciences Cognitives et Adaptatives (LNCA), CNRS, UMR 7364, Strasbourg, France

\* These authors contributed equally to this work.

\* [jyotika.bahuguna@gmail.com](mailto:jyotika.bahuguna@gmail.com) (JB); [timothyv@andrew.cmu.edu](mailto:timothyv@andrew.cmu.edu) (TV); [jonrubin@pitt.edu](mailto:jonrubin@pitt.edu) (JER)



## OPEN ACCESS

**Citation:** Bahuguna J, Verstynen T, Rubin JE (2025) How cortico-basal ganglia-thalamic subnetworks can shift decision policies to increase reward rate. PLoS Comput Biol 21(11): e1013712. <https://doi.org/10.1371/journal.pcbi.1013712>

**Editor:** Alex Roxin, CRM: Centre de Recerca Matemàtica, SPAIN

**Received:** November 12, 2024

**Accepted:** November 6, 2025

**Published:** November 20, 2025

**Copyright:** © 2025 Bahuguna et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data availability statement:** The network codebase utilized in this study can be found on our GitHub repository and accessed at <https://github.com/CoAxLab/CBGTP/blob/main>. Detailed installation

## Abstract

All mammals exhibit flexible decision policies that depend, at least in part, on the cortico-basal ganglia-thalamic (CBGT) pathways. Yet understanding how the complex connectivity, dynamics, and plasticity of CBGT circuits translate into experience-dependent shifts of decision policies represents a longstanding challenge in neuroscience. Here we present the results of a computational approach to address this problem. Specifically, we simulated decisions during the early learning process driven by CBGT circuits under baseline, unrewarded conditions using a spiking neural network, and fit an evidence accumulation model to the resulting behavior. Using canonical correlation analysis, we then replicated the identification of three control ensembles (*responsiveness*, *pliancy* and *choice*) within CBGT circuits, with each of these subnetworks mapping to a specific configuration of the evidence accumulation process. We subsequently simulated learning in a simple two-choice task with one optimal (i.e., rewarded) target and found that, during early stages of learning, feedback-driven dopaminergic plasticity on cortico-striatal synapses effectively increases reward rate over time. The learning-related changes in the decision policy can be decomposed in terms of the contributions of each control ensemble, whose influence is driven by sequential reward prediction errors on individual trials. Our results provide a clear and simple mechanism for how dopaminergic plasticity shifts subnetworks within CBGT circuits so as to increase reward rate by strategically modulating how evidence is used to drive decisions.

instructions and a comprehensive list of implemented functions can be found in the README.txt file within the repository. All datasets generated and analyzed during the course of this research, along with a demonstration demo will be openly available on GitHub at [https://github.com/jyotikab/CBGT\\_maximize\\_RR](https://github.com/jyotikab/CBGT_maximize_RR).

**Funding:** JB is supported by ANR-CPJ-2024DRI00039. TV, JB and JER are partly supported by NIH awards R01DA053014 and R01DA059993 as part of the CRCNS program. JER is partly supported by NIH award R01NS125814, also part of the CRCNS program. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

The task of selecting an action among multiple options can be framed as a process of accumulating streams of evidence, both internal and external, up to a decision threshold. A decision policy can be defined by the unique configuration of factors, such as accumulation rate and threshold height, that determine the dynamics of the evidence accumulation process. In mammals, this process is thought to be regulated by low-dimensional subnetworks, called control ensembles, within the corticobasal ganglia-thalamic (CBGT) pathways. These control ensembles effectively act by tuning specific aspects of evidence accumulation during decision making. Here we use simulations and computational analysis to show that synaptic plasticity at the cortico-striatal synapses, mediated by choice-related reward signals, adjusts CBGT control ensemble activity in a way that improves accuracy and reduces decision time to increase of reward rate during learning.

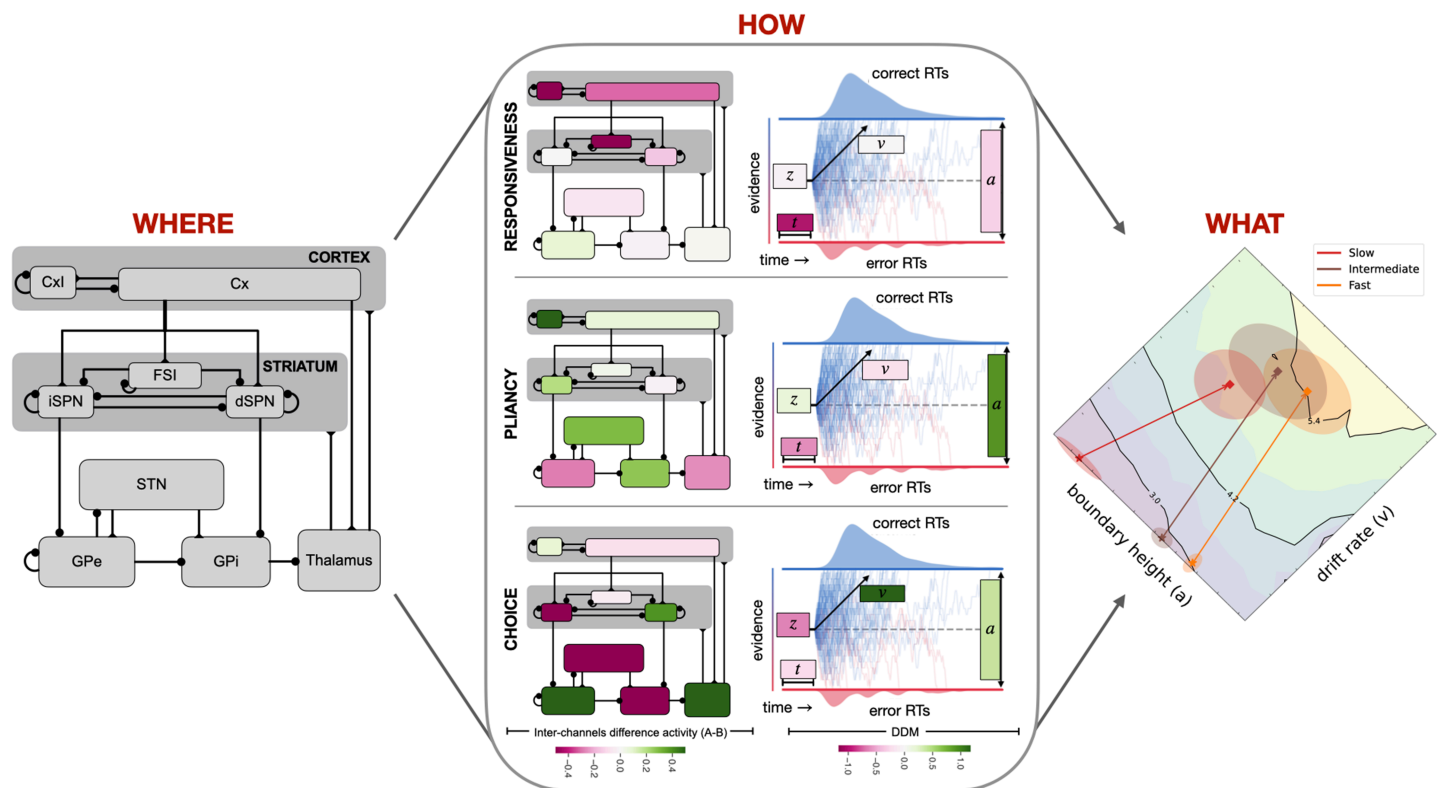
## Introduction

A characteristic of nearly all mammals is the ability to quickly and flexibly shift how currently available evidence is used to drive actions based on past experiences [1]. For example, feedback may be used to shift between making exploratory decisions, where low-value actions are sampled to gain information, and exploitative decisions, where high-value actions are taken to maximize immediate rewards [2–4]. Orthogonal to this exploration-exploitation dimension is a complementary choice about decision speed: actions can be made quickly or slowly depending on immediate goals and confidence level [5]. These shifts between fast or slow and exploratory or exploitative decision policies can be interpreted as different states of an underlying evidence accumulation process [6,7], often captured by mathematical models such as the drift diffusion model (DDM; [8–12]). From this perspective, the values of DDM parameters, such as the drift rate ( $v$ ; the rate of evidence accumulation during a single decision) and boundary height ( $a$ ; the amount of evidence needed to trigger a decision) can be tuned to capture a particular decision policy. Thanks to this mapping, specific ( $a, v$ ) pairs effectively correspond to positions on a manifold of possible decision policies that determine how both internal and external evidence combine to drive eventual actions (Fig 1, “what” panel). Although speed and accuracy are negatively correlated *a priori*, the goal of learning is to converge to a position on this manifold of decision policies that manages to optimize both speed and accuracy for a given task context [13–15].

This form of learning is managed, at least in part, by the cortico-basal ganglia-thalamic (CBGT) circuits, a distributed set of interconnected brain regions that is ideally situated to influence nearly every aspect of decision-making [16–20] (Fig 1, “where” panel). The canonical CBGT circuit includes a collection of interacting basal ganglia pathways that receive cortical inputs and compete for control of an output region (predominantly the internal globus pallidus, GPi, in primates or the substantia nigra pars reticulata, SNr, in rodents) that impacts thalamocortical or superior

collicular activity to influence actions [21–23]. The balance of this competition is thought to map to a configuration of the evidence accumulation process [7,24–29]. Therefore, if behavioral flexibility reflects the *what* and CBGT circuits represent the *where* of flexible decision-making, then we are left with an open question of *how*: how do CBGT circuits control flexibility in decision policies during learning?

In prior work we showed how dynamics of CBGT circuits can be expressed in terms of three subnetworks, or patterns of differential activation across CBGT populations, called *control ensembles*. Each control ensemble tunes specific configurations of the evidence accumulation process, manifested as control over distinct dimensions of a decision policy [29]. In theory, these control ensembles, dubbed *responsiveness*, *pliancy*, and *choice* (Fig 1, “*how*” panel), provide candidate mechanisms for implementing shifts in decision policies during learning. Here we illustrate how a single plasticity mechanism acting at the cortical inputs to the basal ganglia can, through network interactions, leverage the control ensembles to steer behavior during learning during the initial stages of learning. To this end, we simulated a biologically-constrained



**Fig 1. Decision-making deconstructed.** Most voluntary decision policies depend on the CBGT circuits (*where*; left panel). These circuits comprise distributed neuronal populations within the basal ganglia, that interact with each other, as well as cortical and thalamic neurons (connections with circles: inhibition; connections without circles: excitation). This interaction can be described at the algorithmic level by a set of parameters in a process model (here, the DDM) that abstractly simulates evidence accumulation. The goal of this process is to determine the distributions of decision outcomes such as reward rates (*what*; right panel). Contours were generated by simulations of the DDM with drift rate  $v$  and boundary height  $a$  selected on a fine grid of values. Other DDM parameters (onset time,  $t$ ; bias  $z$ ) were fixed. Different initial parameter values and changes in parameters map to different changes in these decision outcomes (arrows in right panel). Control ensembles within CBGT circuits effectively determine the relative configuration of decision policy parameters (*how*; middle panel) [29]; that is, each ensemble represents a mapping between a pattern of increases (green) or decreases (magenta) in firing in CBGT regions (middle panel, left column) and increases (green) or decreases (magenta) of DDM parameters (middle panel, right column). What remains unclear, and we address in this work, is how learning modulates the balance between control ensembles in a way that shifts decision policies so as to maximize reward rate. Cx, cortical PT cells; CxI, inhibitory interneurons; FSI, fast spiking interneurons; d/ISPN, direct/indirect spiny projection neurons; STN, subthalamic nucleus; GPe, external globus pallidus; GPi, internal globus pallidus.

<https://doi.org/10.1371/journal.pcbi.1013712.g001>

spiking CBGT model that learns to select one of two actions via dopamine-dependent plasticity, driven by reward prediction errors across 15 trials, at the cortico-striatal synapses. We then implemented an upwards mapping approach [30], in which the behavioral features (decision times and choices) produced by the simulated CBGT network were modeled across stages of learning using the DDM (see [28,29,31]); specifically, the resulting RT and choice distributions after each stage of learning were used to fit the DDM parameters. Finally, we used various analytical approaches to replicate the existence of the low-dimensional control ensembles prior to learning and quantify how their influence levels change over the course of training. Our results show that value-based learning tunes the influence of CBGT control ensembles to boost reward rate, achieving near-maximal gains across successive decisions when option reward probabilities are well separated.

## Results

### Feedback learning in CBGT networks efficiently increases reward rate

Suppose that an agent encounters a new environment for which it has no relevant prior experience or bias, so that the selection of all options is equally likely at first. In a simple two-choice bandit task, with one rewarded and one unrewarded option, this unbiased starting point would correspond to a 50% error rate. With learning it should be possible to make fewer errors over time, leading to increased rewards, but exactly how this is achieved in practice depends on the decision policy that the agent adopts. For example, if the agent prioritizes speed over all else in its action selection, then its error rate will likely remain high, leading to fewer rewards over time. Conversely, by making sufficiently slow decisions, the agent may be able to achieve an extremely low error rate, leading to greater likelihood of reward on individual trials, but if the response speed is too slow then the rate of reward return over time may take a significant hit. The overall reward rate achieved by the agent thus depends on both decision speed and accuracy. Intuitively, this may be optimized for a fixed level of experience via some compromise between these two dimensions [13,32].

To understand how optimized speed and accuracy emerge from CBGT circuits, we simulated 300 instances of a spiking computational model of CBGT pathways. To generate these instances, we started from parameter ranges defined in our past work [29] (with slight updates - see S2 Table) and used genetic algorithms (see *Materials and Methods – Genetic algorithms*) to derive a collection of networks that operate in a cortico-basal ganglia-thalamic driven regime (i.e., rather than bypassing the basal ganglia to have thalamic activity dictated by its cortical input), that have average firing rates of all relevant cell types within known biological ranges (S1 FigA), and that exhibit a heterogeneous range of response times (S1 FigB). The networks performed a two-armed bandit task with deterministic reward feedback (i.e., the reward probability was 100% for the optimal choice and 0% for the suboptimal one). Using a deterministic reward task, as opposed to a task where rewards are delivered probabilistically, explicitly ties accuracy to reward return and also makes the optimal learning strategy a simple “win-stay/lose-switch” policy [33]. Learning in the network was implemented with dopamine-dependent plasticity at the cortico-striatal synapses, where the magnitude of the phasic dopamine response following each decision was based on reward prediction error (for details see Section: *Dopamine-dependent plasticity of cortico-striatal synaptic weights* and [34]). It should be noted that despite this being a deterministic task, there is ample variability in performance due to noisy background input provided to all populations in the network. Even after 15 trials of learning, by which point learning saturates (S4 Fig), the networks do not reach perfect performance, averaging  $\approx 90\%$  accuracy (S3 Fig).

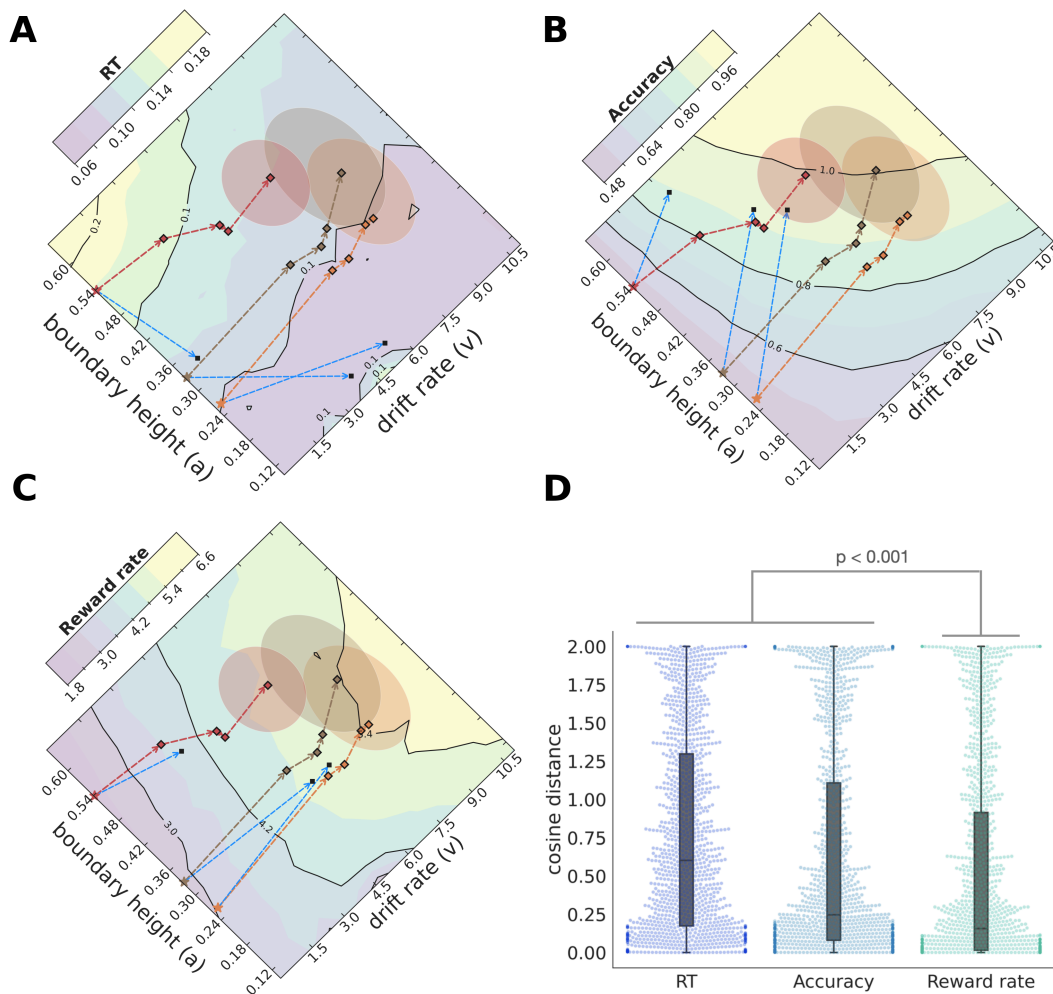
Following a set of simulated trials in this task, we fit the reaction time (RT) and choice probabilities of each network, together reflecting the form of its decision policy, with a hierarchical version of the DDM [35,36]. The DDM provides an intuitive framework for mapping behavioral responses to an evidence-accumulation representation of the decision policy that can be described by only a few parameters [8]. Although there are many possible variants of the DDM that we could use, including versions with collapsing bounds [37–39] or trial-wise evolution of specific parameters [40], we did not include such extensions because our task does not involve factors like urgency or non-stationarity of task states. To try to

capture changes in DDM parameters with learning, we fit the behavioral responses across the 15 learning trials with both the standard static DDM and the reinforcement learning (RL)-DDM variant – specifically the HDDMrl implementation from the HDDM package – which allows for drift rate to change during the learning process based on variations in relevant task or performance parameters. Both approaches yielded poor fits (S6 Fig), likely due to the small number of trials available before learning saturated. Thus, we instead opted for an alternative approach in which, after each predetermined step in learning (2, 4, 6, and 15 trials with plasticity on), we froze the network by turning off plasticity, simulated 300 trials to generate an RT distribution and choice probabilities from the current state of the network, and fit the static DDM and the RL-DDM to these behavioral measures; because performance of the static DDM was substantially better than that of RL-DDM (S6 FigA, S6 FigB), we henceforth focus on the static DDM results. After these probe trials, learning was turned back on and the task progressed. This process captured the behavioral data with high fidelity (S6 Fig) and allowed us to plot each network's performance as a trajectory in the DDM parameter space.

Because of the mapping from network behavioral responses to DDM parameters, we will refer to the 2-dimensional plane of drift rates ( $v$ ) and boundary heights ( $a$ ) as a decision policy manifold. Fig 2 shows the average trajectories of three groups of networks on the ( $v, a$ ) decision policy manifold. For each  $v$  and  $a$  we also estimated the average RT (Fig 2A), accuracy (Fig 2B) and reward rate (Fig 2C; see also [14]). The three network groups shown in this figure represent a tertile-based split of the full set of simulated networks into fast (short RT, orange), intermediate (medium RT, brown), and slow (long RT, red) groups, based on their initial RT values (S1 FigB). We implemented this split to determine whether decision policy adjustments due to learning were influenced by initial biases in the networks. Despite their initial speed differences, all three network groups showed chance level performance before plasticity (S1 FigC) and converged to similar regions of the ( $v, a$ ) space with learning (Fig 2, shaded ellipses). A comparison of behavioral measures and DDM parameters before and after plasticity is presented in S3 Fig.

These trajectories clearly demonstrate that our CBGT network can learn from simple dopaminergic feedback at the cortico-striatal synapses. But what exactly is the objective being optimized by the network? To address this question, we compared the change at each step of learning to the predicted direction that the network would take if it were maximizing one of the three possible behavioral objectives: speed, accuracy, or reward rate. Note that we can plot contours for each of these quantities (using RT as a gauge of speed) over the ( $v, a$ ) domain. Although the mapping between ( $v, a$ ) and either speed, accuracy, or reward rate is not bijective, we will nonetheless refer to the ( $v, a$ ) plane as the speed manifold, accuracy manifold, or reward rate manifold when it is shown along with the contours of the corresponding quantity. The predicted directions of objective maximization are illustrated as blue vectors in Fig 2A–2C, reflecting steps from each initial point that are in the direction of the gradient of each objective (i.e., the direction of maximal change, which lies orthogonal to the contours, shown with the same length as the vector representing the actual network evolution at the first step of learning in each case). Analysis of the trajectories in Fig 2A reveals that while plasticity decreases RTs with learning, the angles of the learning trajectories do not align with the optimal directions for maximally reducing RT. Similarly, the network trajectories do not align with the vectors that would be expected if they were maximizing accuracy alone (Fig 2B). In contrast, the average trajectories along the reward rate manifold were closest to the gradient and hence to the optimal trajectories attainable for that manifold (Fig 2C). Moreover, the rate of increase in reward rate was similar regardless of the network's initial speed bias.

To quantify the alignment of observed network trajectories to the expected directions of maximal change, we calculated the cosine distance between the observed vector and the optimal vector, normalized to the observed vector's length, at each learning step. While there is substantial variability across networks (Fig 2D), there was a consistent effect of objective type on the network fits ( $F[3597, 2]=47.2$ ,  $p < 0.0001$ ). Fits to the reward rate trajectories were consistently better than to either speed ( $t(299)=13.59$ ,  $p < 0.0001$ ) or accuracy ( $t(299)=8.35$ ,  $p < 0.0001$ ) trajectories. This effect held regardless of a network's initial speed (S4 Fig). This effect also held when the network was trained on non-deterministic reward probabilities (S5 Fig), except for high reward uncertainty such as 75%, where reward rate and accuracy fits approximately equalized (see Discussion).



**Fig 2. Dopamine-dependent cortico-striatal plasticity drives CBGT networks in the direction of reward rate maximization.** (A) The evolution of RTs achieved by a DDM fit to CBGT network behavior, projected to  $(v, a)$ -space. The average starting position for the fast (orange), intermediate (brown) and slow (red) networks are shown as stars. The squares indicate the evolution of each network group over the plasticity stages, which converge after 15 trials (shaded elliptical regions). The yellow (purple) colors represent high (low) RTs. The network trajectories do not evolve in the direction that would be expected to minimize the RTs (e.g., optimal direction shown in blue from the initial position of all three speed groups). (B) The yellow (purple) colors represent high (low) accuracy. The networks evolve towards increasing expected accuracy but not in an optimal fashion (trajectories vs. blue arrows). (C) The yellow (purple) colors represent high (low) reward rate. The network evolution aligns closely with the direction that maximizes the reward rate (blue arrows). (D) The cosine distances calculated for every network at each plasticity stage for RT, accuracy and reward rate were pooled together and shown as distributions.

<https://doi.org/10.1371/journal.pcbi.1013712.g002>

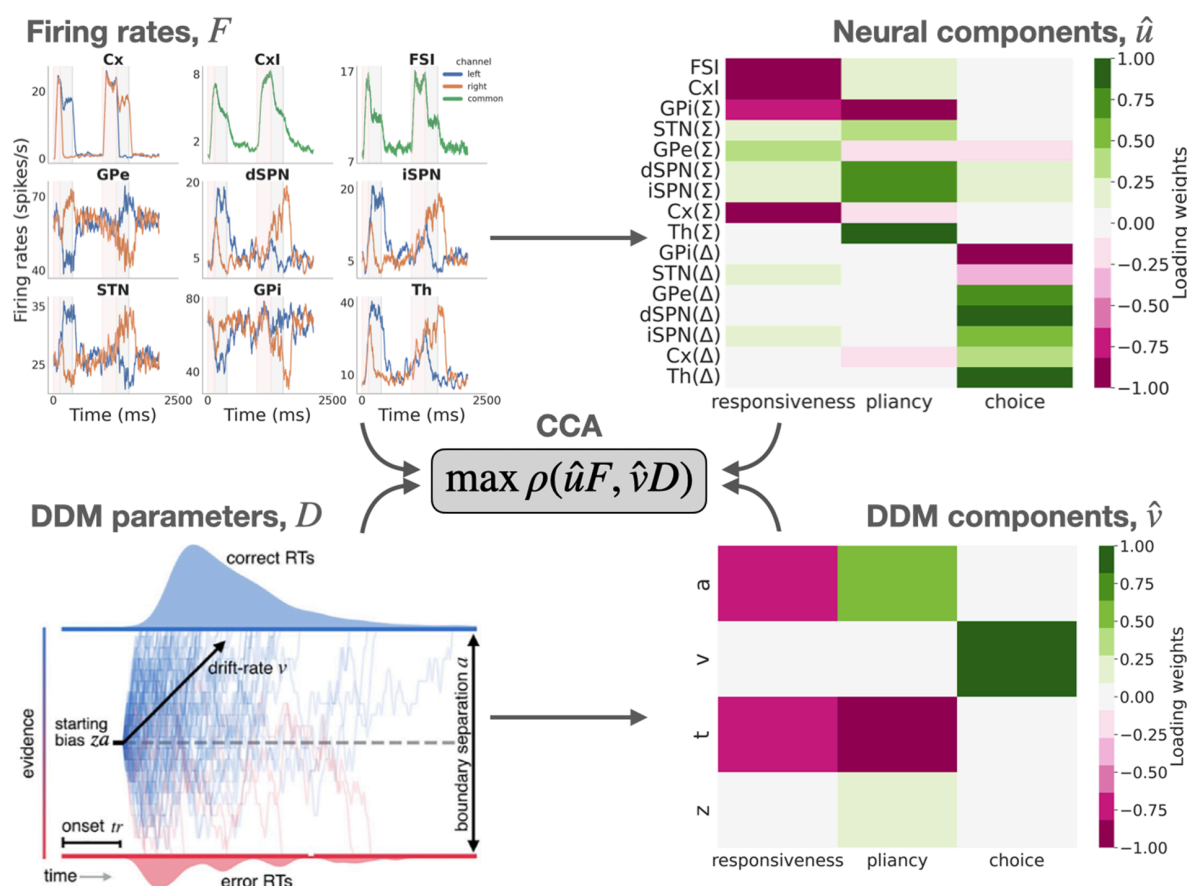
Overall, our DDM results show that our biologically detailed model of the CBGT circuit can effectively learn to increase reward rate by managing the accuracy and speed of the evidence accumulation process via dopaminergic plasticity at the cortico-striatal synapses. Moreover, this increase is near-maximal when reward probabilities are deterministic.

### Low-dimensional control ensembles that map to general decision policies

The CBGT network and DDM are, respectively, implementation-level and algorithmic-level descriptions of the evidence accumulation process that guides goal-directed behavior. We have previously shown that there is a low-dimensional, multivariate mapping between these two levels of analysis in the absence of learning [29]. Here we set out to replicate

this observation with the CBGT parameter sets used in the current study, with the aim of analyzing their contributions to the dopaminergic learning process. For this step, we considered two aspects of activity within each CBGT population: global activation across the two action representations (sum of the activity in that region, across both channels;  $\Sigma$ ) and bias towards one action representation (difference in activity within each region, across the action channels;  $\Delta$ ). Using canonical correlation analysis (CCA), we captured the low-dimensional components that maximally correlate variation in CBGT activity with variation in DDM parameters. This analysis identified three such components (Fig 3). We refer to these low-dimensional components as *control ensembles*.

The three control ensembles identified by our analysis nearly perfectly replicate our prior work [29], where they are described in more detail (see also Section *Upward mapping*). Thus we kept the labels *responsiveness*, *pliancy*, and *choice* ensembles for the first, second, and third components, respectively. The recovered components are shown in both CBGT and DDM parameter spaces in Fig 3 (right panels). The responsiveness component describes the agent's sensitivity to evidence, both in terms of the delay before the agent starts to accumulate evidence ( $t$ ) and how significantly the presence of evidence contributes to achieving the decision threshold ( $a$ ). The dominant features of CBGT activity that vary along the responsiveness control ensemble loadings are a global inhibitory signal, including fast-spiking interneuron (FSI)



**Fig 3. Canonical correlation analysis (CCA) identifies control ensembles (cf. [29]).** Given matrices of average firing rates,  $F$  (both summed rates across channels,  $\Sigma$ , and between-channel differences,  $\Delta$ ), and fit DDM parameters,  $D$ , derived from a set of networks at baseline (left panels), CCA finds the low-dimensional projections,  $\hat{u}$  for firing rates and  $\hat{v}$  for DDM parameters (right panels), which maximize the correlation,  $\rho$ , between the projections  $\hat{u}F$  and  $\hat{v}D$  of  $F$  and  $D$ . Blue lines in the  $F$  plot show left channel activity, orange show right channel activity, and green shows populations that go across both channels.

<https://doi.org/10.1371/journal.pcbi.1013712.g003>

and overall internal globus pallidus ( $GPI(\Sigma)$ ) activity, as well as overall excitatory and inhibitory cortical activity ( $Cx(\Sigma)$ ,  $CxI$ ). Because the dominant CBGT and DDM loadings for the responsiveness control ensemble have the same sign (all negative), they imply that a *decrease* in the weighted activity of the loaded populations corresponds to an *decrease* in onset time,  $t$ , and  $a$  and, hence, to an *increase* in overall responsiveness.

The pliancy component refers to the level of evidence that must be accumulated before committing to a decision. As with responsiveness, pliancy loads mostly on  $a$  and  $t$ , but now with opposing signs for these two loadings, corresponding to the idea that even though an agent is attentive to evidence (small  $t$ ), it requires a substantial accumulation of evidence to reach its threshold (large  $a$ ). The CBGT activity features that characterize pliancy are the overall engagement of the BG input nodes (i.e., global dSPN and iSPN activity, with a smaller STN contribution), as well as total GPI and thalamic activity, with oppositely signed loadings to each other. For the pliancy component, a change in the activity consistent with the cell type loadings (e.g., increase in SPN activity) corresponds to a decrease in overall pliancy (e.g., increase in  $a$ ).

Lastly, the choice component represents the intensity of the choice preference and is reflected largely in  $v$  and the neural correlates of competing choice representations in the CBGT (i.e., differences in activity across the two action channels within each BG region). A change in activity consistent with the cell type loadings (e.g., greater difference in dSPN activity between the two channels) corresponds to a stronger commitment towards the more rewarded option (i.e., larger  $v$ ).

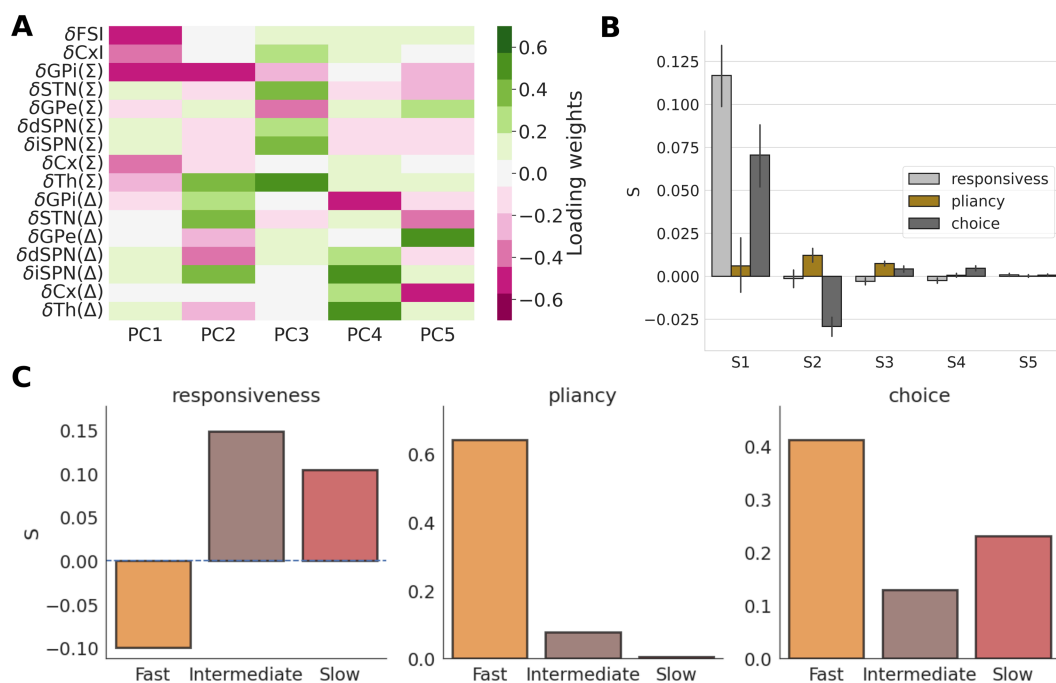
In summary, each CBGT control ensemble can be interpreted as specifying a coordinated collection of changes in CBGT neural activity that modulates specific aspects of the decision policy, as represented by the DDM parameters needed to fit that behavior. Now that we have delineated the control ensembles embedded within the CBGT network (cf. [29]), we are ready to consider how dopamine-dependent plasticity regulates their influence in a way that collectively drives decision policies to maximally increase reward rate.

### Cortico-striatal plasticity drives control ensembles during learning

Our analysis of the CBGT network behavior (Fig 2) shows that dopamine signaling at the cortico-striatal synapses is enough to elicit changes in the evidence accumulation process that increases reward rate. This observation suggests that there are emergent driver mechanisms, originating from cortico-striatal synaptic changes, that tune the control ensembles in a way that achieves this outcome. That is, if each control ensemble represents a knob to tune an aspect of the decision policy, then a driver mechanism selects a set of adjustments of the knobs that yields an overall decision policy selection. We next set out to identify these emergent drivers.

As a first step toward quantifying the modulation of CBGT activity after 15 learning trials, we calculated the principal components of the overall change in firing rates of all 300 networks. The first 5 of these components collectively explain more than 90% of the observed variance (S7 FigA, thick blue line marked "All"). The loading weights (Fig 4A) show that the first and third components reflect the global activity of subsets of CBGT nuclei. The second, fourth and fifth components relate more strongly to the bias towards one option, with predominant loadings on differences in rates across action channels. Together, these components represent the collection of changes in firing rates that result from learning-related changes at the cortico-striatal synapses.

We next calculated the matrix  $S$  of weighting factors (*drivers*) for the firing rate components, which describe what combination of adjustments to the control ensembles best account for the associated firing rate changes (Fig 4B; for full description of this approach, see Methods subsection *Modulation of control ensembles by plasticity*). To interpret the drivers of control ensemble influence (Fig 4B), it is important to note that positive (negative) coefficients correspond to changes in control ensemble activity in the same (opposite) direction as indicated by the loadings in Fig 3. The first driver corresponds to a large amplification of the responsiveness control ensemble, and hence a decrease in various forms of global inhibition in the CBGT network. The first driver is also associated with a boost to the choice control ensemble, increasing the bias towards the rewarded choice. The second driver has a strong negative weight on the choice control ensemble and a positive weight on the pliancy control ensemble. The third, fourth and fifth drivers feature weaker effects,



**Fig 4. Plasticity-induced changes of control ensemble influence.** (A) The loading weights of the first 5 PCs of firing rate changes from before to after plasticity, pooled for all networks. (B) The drivers (columns of  $S$ ), which quantify the modulations of control ensembles (responsiveness, pliancy, choice) that capture each PC (pooled for all network classes). (C): The variance-weighted drivers for the three control ensembles, computed separately for the three network classes (fast, intermediate and slow).

<https://doi.org/10.1371/journal.pcbi.1013712.g004>

with small modulations of all three control ensembles. Based on this analysis across all of the networks, the overall modulation of the control ensembles due to plasticity, calculated as the weighted sum over all drivers, adjusted by the % of variance explained by each PC, is shown in S7 FigB. All three control ensembles end up being boosted. This means that, to varying extents, the activity measures that comprise these ensembles change in the directions indicated by their loadings in Fig 3. In this way the general trend is for the CBGT networks to become more responsive, yet less pliant, which together amount to an earlier onset of evidence accumulation without much change in boundary height. Coincident with this, we also see that the CBGT networks exhibit more of an emergent choice bias with learning.

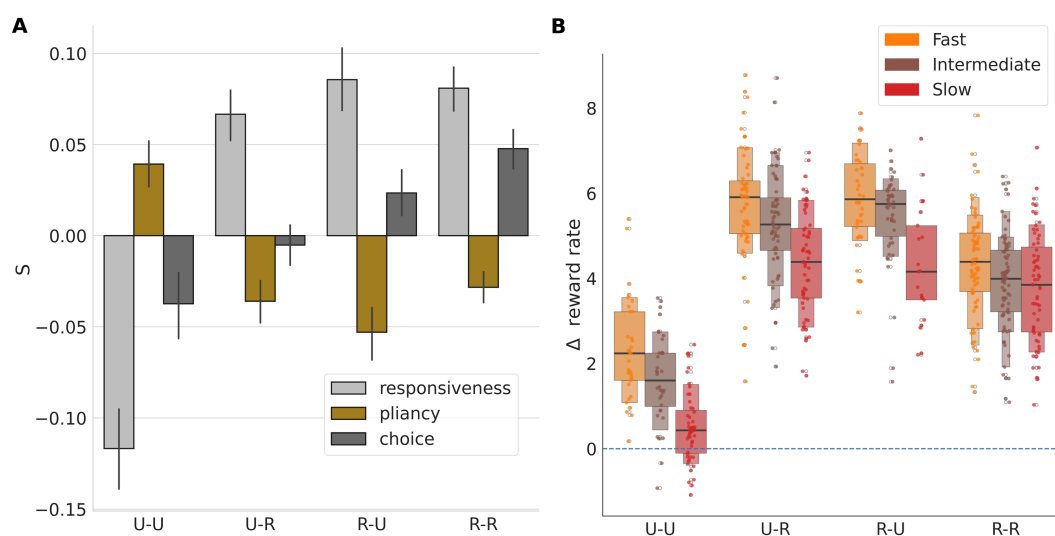
Because of the difference in decision policies across the fast, intermediate, and slow networks, we recomputed the drivers separately for each network type. This was done by considering the firing rate differences ( $\Delta F$ ) and calculating the  $S$  loadings for fast, intermediate, and slow networks separately (see Methods - section *Modulation of control ensembles by plasticity*). The explained variance for each of the three network types is shown in S7 FigA, and their corresponding PCs and goodness of fits are shown in S8 Fig. As expected, the drivers show different changes across the network types (Fig 4C). The driving factor corresponding to responsiveness is negative for fast networks, while remaining positive for the others. The pliancy and choice factors were positive for all three networks, but pliancy was by far the largest for fast networks and quite small for the other two network types. Referring to the DDM parameter changes associated with changes in control ensemble loadings (Fig 3), we see that the decrease in loading of responsiveness and strong increase in loading of pliancy for fast networks would both promote an increase in boundary height,  $a$ . This aligns with the fact that, of the three network types, only fast networks show an increase in  $a$  over the course of learning (Figs 2 and S9 Fig). The DDM parameters exhibit saturation after approximately 15 learning trials, as shown in S9 FigA–S9 FigC, as a result of the maximum bound imposed on the corticostriatal weights, shown in S3 Table. The upper limit on the corticostriatal weights is

the result of the multiplicative formulation of weight update (Eq 2 in section *Dopamine dependent plasticity of corticostriatal weights*). This ensures that the weight updates are smaller as the weights approach certain maximum or minimum bounds. Multiplicative update rules are widely used in plasticity models and have been shown to produce unimodal weight distributions and stable learning dynamics [41,42]. Overall, we see that the specific way that plasticity adjusts the weighting of the control ensembles to drive changes in decision policies depends on the initial tuning of the network. Since plasticity results from the sequence of decisions and rewards that occur during learning, we next investigate more directly how specific decision outcomes lead to this dependency.

### The influence of feedback sequences on driving of control ensembles

In the previous section we described the overall effects of cortico-striatal plasticity on control ensemble tuning. We now turn to analyzing the early temporal evolution of these effects by focusing on the initial two learning trials. Specifically, we examined the modulation of the control ensembles for different combinations of successes (i.e., rewarded trials; R) and failures (i.e., unrewarded trials; U) achieved by the first two consecutive choices. For this analysis, we implemented our usual DDM fitting process followed by CCA for networks that were frozen (i.e., with plasticity switched off) after two trials, and we grouped the results based on the sequence of choice outcomes. The drivers for each sequence of outcomes, U-U, U-R, R-U and R-R, are the linear combinations of columns of  $S$  that express the corresponding changes in firing rates in terms of adjustments to the control ensembles (Fig 5A).

First, consider the case of networks that receive no rewards (U-U). Here we infer that the boundary height increases, due to a simultaneous decrease in driving of the responsiveness ensemble and increase in driving of the pliancy ensemble, both of which result in a boost of the boundary height. In addition, driving of the choice ensemble is reduced. Thus, two consecutive unsuccessful trials yield an overall increase in the degree of evidence needed to make a subsequent decision by simultaneously increasing the boundary height and decreasing the drift rate. Moreover, slow networks encounter U-U outcomes more often than other network classes in the first two trials (S1 Table), which presumably constrains the increase in responsiveness and choice seen in these networks during learning (Fig 4C). On average, however,



**Fig 5. Suboptimal and optimal choices modulate control ensembles in opposite directions.** (A) The modulation of control ensembles associated with various reward sequences encountered in two initial trials with cortico-striatal plasticity. U represents “Unrewarded” and R represents “Rewarded” trials. (B) The reward rate changes obtained by simulation of networks with synaptic weights frozen after various reward sequences occurred on two initial trials.

<https://doi.org/10.1371/journal.pcbi.1013712.g005>

fast networks make more mistakes than the other networks. This result, which we can display graphically in terms of the proportion of unrewarded trials, or mistakes, encountered after the first two plasticity trials (S9 FigD), likely explains the negative loading for responsiveness and high positive loading for pliancy for fast networks shown in Fig 4C.

In contrast, two consecutive successful trials (R-R, far right of Fig 5A) produce largely the opposite effect. The influences of the responsiveness and choice ensembles increase, resulting in lower onset time and boundary height, along with an increase in the drift rate. This coincides with a weak change in pliancy. As a result, in the R-R case, the decision policy is tuned to include a decreased degree of evidence needed to make subsequent decisions.

Not surprisingly, the two mixed combinations of outcomes (U-R, R-U) have largely similar effects on the responsiveness and pliancy ensembles, regardless of the order of outcomes. In both cases the loading of responsiveness increases and that of pliancy decreases, resulting in less overall evidence needed to trigger a decision (by shrinking the boundary height, without much change in the onset time). However, when the first trial is unsuccessful (U-R) the influence of the choice ensemble decreases, while it increases when the first trial is successful (R-U). Indeed, looking at the progressive change in the choice ensemble across the four unique sequences of trials, it appears that early success (i.e., reward in the first trial) boosts the choice ensemble influence while early failure (i.e., unrewarded first trial) does the opposite. When these combined drivers are recomputed separately for each network class, the learning-induced modulations of the ensembles follow the same general trend (S10 Fig), with quantitative details depending on the network class.

The preceding analysis shows how the relative contributions of the control ensembles to the evidence accumulation process depend on trial outcomes. What are the results of these changes on the performance of the network? To illustrate these effects, we plot the distribution of changes in reward rates associated with each set of outcomes and separate by network types in Fig 5B. Although all distributions are generally positive, there is significant variation in reward rate changes across the different feedback sequences ( $F(586, 3) = 254.4, p < 0.0001$ ). The reward rate also varies significantly with the network type ( $F(586, 2) = 46.8, p < 0.0001$ ), and the interaction term between network types and feedback sequences is significant as well ( $F(586, 6) = 3.8, p = 0.001$ ). Compared to all other conditions, the networks that made two consecutive unsuccessful choices (U-U) yielded the smallest changes in reward rates (values of all network types pooled together, all two-sample  $t(319) < -18.27$ , all  $p < 0.0001$ ). The two mixed feedback conditions (U-R, R-U) had higher growth in reward rates than the condition with two rewarded trials (R-R; all  $t(384) > 8.1$ , all  $p < 0.001$ ), because mixed conditions not only lead to strengthening of the correct choice but also weakening of the incorrect choice, unlike R-R which only leads to the former. In all cases, the trend was for faster networks to achieve greater increases in reward rate than slower networks. As expected, the impact of feedback sequences on reward rate is associated with underlying changes in both accuracy (S11 FigA) and decision speed (S11 FigB). Like reward rate, the increase in accuracy was highest for the mixed feedback conditions (U-R, R-U) due to the combined strengthening of the correct choice and weakening of the incorrect choice. Two consecutive unsuccessful choices (U-U) represents the only condition that leads to an increase in decision times, expressed as negative  $\Delta RT$ s. This outcome is consistent with the increase in boundary height that occurs in this case, whereas all other feedback conditions lead to a decrease in decision times.

## Discussion

Adaptive behavior depends on flexible decision policies (*what*), driven by CBGT networks (*where*) that shift their activity in order to increase reward rate by coordinated adjustments of a set of underlying control ensembles (*how*; Fig 1). In this work, we focused on the *how* part of this process, employing a mapping upward in abstraction between a biologically realistic model of CBGT pathways and the DDM. This approach helps to reveal the complex, low-dimensional structure of CBGT subnetworks that influence decision-making policies (Fig 3). Specifically, we recapitulated recent results [29] showing the existence of three main CBGT control ensembles shaping decision-making that represent *responsiveness*, *pliancy*, and *choice* (direct vs. indirect pathway competition; Fig 3) and serve to regulate the process of converting evidence accumulation into action selection. We then showed how, within our model, driver mechanisms tune these control ensembles

strategically during early learning (Figs 4 and 5) in order to increase reward rate. Moreover, although they all optimize the same quantity (reward rate), we found that modulation of control ensembles differs across networks depending on their *a priori* decision policy (fast, intermediate, or slow). While plasticity increases responsiveness and choice in all networks, to varying extents, fast networks alone decrease responsiveness (Fig 4C) and correspondingly increase boundary height (a; S9 FigA). Put together, our results propose a new framework for understanding how subnetworks within CBGT circuits can dynamically regulate decision-making, driven by dopaminergic plasticity at the cortico-striatal synapses.

Perhaps the most surprising aspect of this theoretical analysis is the sophisticated adjustments that emerged from a simple plasticity mechanism acting on just one class of CBGT synapses. Dopaminergic learning at the cortico-striatal synapses was sufficient to push our naive networks from an exploratory decision policy to an exploitative policy. This progression did not maximize the instantaneous reward rate as discussed in [13,15]. Instead, it jointly resolved speed and accuracy in a way that aligned with the reward rate gradient, corresponding to changes in reward rate along the direction of optimal increase (Fig 2). This pattern is consistent with findings from perceptual learning tasks in rats [32]. Although improvements in accuracy combined with reductions in reaction time naturally lead to higher reward rates, the stronger alignment of the DDM parameter change vector with the reward gradient, rather than with the accuracy gradient, is not simply a by-product of reward rate being defined as the ratio of accuracy to RT.

The reduction in RTs observed during learning was not built into our network *a priori*; rather, it emerged from cortico-striatal plasticity driven by decision accuracy. While classic work has established that learning in various forms of RL algorithms follows the gradient of expected *reward* maximization [43–45], it is not readily apparent that this process also aligns with the gradient of *reward rate* improvement. How, then, can reward signals based only on accuracy lead to an optimal increase of reward rate? The answer lies in the architecture of the CBGT circuits themselves. In our model, synaptic plasticity is restricted to cortico-striatal synapses and follows a mechanism that closely resembles Herrnstein's matching law, where the probability of choosing an option is proportional to the probability of reward derived from selecting that action in the past [46]. However, due to the specific pattern of synaptic coupling across interconnected populations within the CBGT network, changes in activity in these plastic synapses propagate throughout the entire circuit in a structured way. A key emergent property of this architecture, as revealed in our simulated experiments, is a reduction in decision times, even though speed itself is not explicitly rewarded. Thus, our model tends to make slower decisions early in learning, but improves accuracy and becomes faster naturally as learning progresses. This progression is similar to behavioral observations in rodents [32], non-human primates [47], and humans [48,49]. Finally, while our findings suggest that this complex behavior arises naturally from dopamine-dependent plasticity at the cortico-striatal synapses together with the intrinsic architecture of the CBGT circuit, the expression of this mechanism may vary depending on task context, including factors such as differences in effort, which we do not model here, or the nature of performance feedback [14,50].

We note that reward rate maximization via cortico-striatal plasticity occurs under the constraints of a lower bound on the environmental and reward uncertainty. With a high enough uncertainty in the reward probabilities (e.g., 75%), we find that the cosine distances for maximizing accuracy and reward rate become comparable (S5 FigB and S5 FigC). In our prior work with humans [6], this reward probability is near the lower bound of detectability of the optimal choice before participants begin random decisions. Like our current simulation experiments, this prior empirical work afforded participants less than 2 dozen trials to learn the reward contingencies (before they were switched). Thus, the learning we are capturing here is the initial process, when the rate of learning is accelerating the fastest. With prolonged training, on the order of hundreds of trials or more, both humans [51,52] and rodents [53] can learn to discriminate more conflicting reward probabilities (e.g., 65% vs. 35%). Within the initial learning window that we explored here, even when the distinction between the optimality of reward and accuracy gains disappears at higher levels of stochasticity, the overall small cosine distance between the direction of reward rate increase and the gradient (S5 FigB) remains a desirable outcome. The relative decline in optimality of the increase in reward rates that we observe in this case comes from slowed reaction times during early phases of learning (as shown in rodents by [32]) and when mistakes are made (as shown in our

simulations - S11 FigB) as a result of an increased decision threshold, which ensures an accumulation of sufficient evidence before commitment to a decision under high reward uncertainty.

Related predictions at the abstract level have been made by models that directly combine reinforcement learning with evolution of the DDM parameters [54,55]. These studies demonstrate that the drift rate depends on the difference in values between optimal and suboptimal actions, which increases with learning. In contrast, the boundary height is proportional to the effective values of the choices and typically shows a slight decrease as learning progresses. Another class of promising models are the reinforcement learning and racing diffusion (RL-RD) models, which can represent multi-choice decision-making with DDM-like accumulation. An evaluation of the efficacy of RL-DDM (*HDDMrl* in the package *hddm*) in capturing the process of plasticity shows that trial-wise modulation of the drift rate is not sufficient to explain the plasticity process in the CBGT networks (S6 Fig). However, some of the other models [56] share a conceptual similarity with our CCA components in that they include q-values related to sum ( $\Sigma$ ) and difference ( $\Delta$ ) elements and may be better suited than RL-DDM to model the plasticity mechanism in CBGT networks. The RL-RD class of models offers alternative options for parameterizing learning data from our CBGT circuit model. For the current work, however, we limited our analysis to estimating DDM parameters across learning stages to maintain consistency with our previous findings [29] and our current results on control ensembles in naïve CBGT networks.

Our primary goal with the analyses described in this paper was to decompose the circuit-level effects of plasticity that underlie adaptive reward rate maximization in terms of learning-related changes in the driving of the control ensembles. Based on the relation of the control ensemble loading to the evidence accumulation parameters (Figs 3 and 4C), the effective learning-related changes result in shorter decision onset delays, higher rates of evidence accumulation, and variable changes in decision threshold as learning progresses (S9 Fig). On the shorter timescale of consecutive trials, each possible set of pairs of reward outcomes induces a specific adjustment of control ensembles in a way that increases subsequent accuracy and reward rate (Figs 5 and S11 Fig). Interestingly, but perhaps not surprisingly, our results predict that mixed feedback, such as one rewarded and one unrewarded trial, will result in a stronger increase in reward rate than two consecutive rewarded trials. This finding is consistent with past results, as well as general intuition, on the benefits of exploration for effective learning [57,58].

It is, however, important to note that cortico-striatal plasticity may explain only a part of the decrease in decision speed seen in experiments. Additional boosts in speed may result from an agent's increased confidence in the outcomes of its decision derived from other information sources [59]. Moreover, an experimental paradigm that requires learning an explicit minimization of decision times may reveal other novel CBGT control ensembles, apart from those that we report here.

The existence of a small set of CBGT control ensembles, and the details of their components, represent some of the key predictions that emerge from our modeling study. Directly recovering these ensembles in real CBGT circuits would necessitate simultaneous *in vivo* recording of nine distinct cell populations across at least five distinct brain regions during a learning task. This is currently outside the scope of available empirical technology. While we hope that future experiments will test more focused aspects of our predictions, we can already extract relevant findings from the extant literature. For example, the predominant loadings in the responsiveness ensemble reflect the level of engagement primarily of input-level (cortical and FSI) components and inhibitory outputs (GPI) of the network, with higher loadings corresponding to less activation (Fig 3). The increase in responsiveness associated with learning in intermediate and slow networks in our model therefore matches the suppression of activity in the subpopulation of striatal FSIs that was observed after learning in non-human primates [60]. Interestingly, experiments have also found evidence for an earlier onset of activity in the striatum with the progression of learning in non-human primates [61]. This is consistent with the decrease in onset time that arises via the learning-induced increase in the responsiveness and pliancy ensembles in all network classes in our simulations.

The pliancy ensemble, reflecting the influence of global striatal activity, including thalamic inputs to the striatum, as well as the influence of STN activity, is associated with the onset time and boundary height parameters; however, unlike the

responsiveness ensemble, the pliancy ensemble has opposing loadings between onset time and boundary height. Thus, an increase in activity of the pliancy ensemble corresponds to an earlier onset of evidence accumulation, but with more information required to trigger a decision. This places an emphasis not on the collection of evidence itself, but instead on the agent's willingness to be convinced by this evidence. It has been shown that an increase in the conflict between action values is associated with an increase in global STN activity [62–64], consistent with a strengthened driving of our pliancy ensemble (see Fig 3, which shows the positive loading for summed STN activity in the pliancy ensemble, and Fig 4A and 4B, showing that an increase in the drive of pliancy in S3 aligns with an increase in STN firing in PC3). Also, because our simulations show an increase in efficacy of the pliancy ensemble with value-based learning for fast and intermediate networks (Fig 4B and 4C), we predict based on the loadings in Fig 3 that the overall level of striatal SPN activity will increase as learning progresses. In contrast, activity in the GPi would decrease. The predominant contributions of this effect are predicted to occur in response to unrewarded trials (Fig 5A). Consistent with these predictions, past studies have shown increases in striatal activity with learning [65]. Related findings have been interpreted as being potentially linked to increased task attentiveness [66] or increased motivation [67,68]. Both effects are consistent with the lowering of onset time associated with our pliancy ensemble. Interestingly, increases in striatal activity, as measured via fMRI, have been found to be beneficial for learning in adolescents [69]; our results suggest that such increases in the pliancy ensemble loading could relate to learning from mistakes (Fig 5A, U-U case).

Finally, the choice ensemble, which corresponds to the degree of competition between direct and indirect pathways across action channels, is strongly associated with drift rate. Consistent with this relationship, single unit activity in dorsal striatum has been shown to reflect the rate of evidence accumulation and consequently preference for a specific response to a stimulus [70]. At the macroscopic scale, we recently found that the competition between action representations in CBGT circuits, measured with fMRI, is indeed reflected in the drift rate in humans [7]. At the causal level, a recent study with patients suffering from dystonia showed that deep brain stimulation in the GPi increased the likelihood of exploratory behavior, which was encoded as decrease in the drift rate [19]. Whether deep brain stimulation (DBS) increases or decreases the output of its target area remains controversial [71–73]. However, based on the loadings in the choice ensemble, we predict that the observed decrease in drift rate corresponds to increased similarity in activity across GPi neurons in different channels—a likely outcome if DBS similarly impacts all channels.

As learning occurs in our model CBGT network, the control ensemble loadings appear to co-evolve. The merit of the control ensemble idea is that it lets us decompose a complicated evolution process into interpretable components. Nonetheless, it can also be informative to consider combined effects that result from the simultaneity of changes in control ensemble loadings. As one example, we note that in non-human primates, stimulation of the caudate nucleus in the striatum reveals a negative correlation between drift rate and boundary height [74]. Our model captures this early negative correlation in learning, where pairs of unrewarded trials decrease the loading of responsiveness and choice while increasing the loading of pliancy (Fig 5A). This shift reflects overall striatal engagement across both channels, potentially mirroring the effects of stimulation, resulting in an increase in boundary height and a decrease in drift rate. In contrast, other outcome pairs produce the opposite effects. Over the longer course of learning in slow and intermediate networks, we observe an increase in drift rate and a decrease in boundary height, with responsiveness and choice ensembles playing a more prominent role. This trend suggests a gradual shift in importance from pliancy (e.g., overall striatal engagement) to responsiveness as learning progresses.

Overall, our results suggest how the low-dimensional substructure of CBGT circuits may adapt behavior during learning by adjusting specific aspects of the evidence accumulation process, thereby influencing the current state of a decision policy. Notably, we demonstrate that dopamine-dependent synaptic plasticity at cortico-striatal synapses, driven by choice-related reward signals, can strategically coordinate control ensemble activity to improve accuracy while reducing decision times, thereby increasing reward rate. As we have discussed, these findings not only align with previous empirical observations but also offer clear predictions for future experimental investigations.

## Materials and methods

### CBGT network

The CBGT model used in this work is a biologically constrained spiking network including neuronal populations from the striatum (dSPNs, iSPNs and FSIs), globus pallidus external segment (GPe), subthalamic nucleus (STN), globus pallidus internal segment (GPi), thalamus and cortex (excitatory and inhibitory components). For a two-choice task, each choice representation is implemented as a “channel” [23,28,29,34,75], so the model includes two populations of each neuron type except FSIs and inhibitory cortical neurons, which are shared across channels. On each trial, the two excitatory cortical populations receive excitatory synaptic inputs, representing evidence related to the available options, from a stochastic spike generator. This process has a baseline rate sampled from a normal distribution with a mean and standard deviation of 2.5 Hz and 0.06, respectively. To this baseline, we add a ramping component representing the presence of some stimulus or internal process that drives the consideration of possible choices. This component rises linearly until it reaches a maximum value ( $f_{target} = 1.0$ ), which was kept constant for all simulations in order to appropriately compare decision times. Specifically, we take

$$f_{ramp}(t) = f_{ramp}(t - dt) + 0.1[f_{target} - f_{ramp}(t - dt)]$$

where  $dt$  is the integrator time step, such that the total frequency of inputs to the cortical populations evolves according to

$$f_{ext}(t) = f_{baseline} + f_{ramp}(t).$$

In all of our simulations, evidence for the two options, as represented by this frequency of inputs to the two cortical populations, was equally strong, such that changes in outcomes across conditions resulted entirely from learning downstream from the cortex.

All populations have an external current  $I_{ext}$  to tune their baseline firing rate, given by

$$I_{ext}(t) = S_{ext,AMPA}(V(t) - V_E) + S_{ext,GABA}(V(t) - V_I)$$

where  $S_{ext,x}$  for  $x \in \{AMPA, GABA\}$  is a mean-reverting random walk derived from the stochastic differential equation

$$dS_{ext,x} = \frac{(\mu_{ext,x} - S_{ext,x})}{\tau_x} dt + \sigma_{ext,x} \sqrt{\frac{2}{\tau_x}} dW_t.$$

Here,  $W_t$  is a Wiener process,  $\tau_x$  is the time decay of the external current, and  $\mu_{ext,x}$  and  $\sigma_{ext,x}$  are computed as

$$\begin{aligned} \mu_{ext,x} &= 0.001 E_{ext,x} f_{ext,x} N_{ext,x} \tau_x, \\ \sigma_{ext,x} &= E_{ext,x} \sqrt{0.0005 f_{ext,x} N_{ext,x} \tau_x}. \end{aligned}$$

The parameter  $f_{ext,x}$  is the external input frequency,  $E_{ext,x}$  is the mean efficacy of the external connections,  $N_{ext,x}$  is the number of connections, and  $\tau_x$  is the time decay constant. Values of all of these parameters are specified in S4 Table.

Specifically, the cortico-striatal projections to both dSPNs and iSPNs in the model were plastic and were modulated by a dopamine-dependent spike timing dependent plasticity rule [31,76,77]. On a trial, a choice was selected if the firing rate in the thalamic population within its action channel reached 30 Hz before the rate of the other thalamic population hit that level. The complete details of this network can be found in our methods paper [34].

## Characterization of networks before plasticity

In our previous work, we identified control ensembles based on extensive simulation of the CBGT network with each of 300 parameter sets selected using Latin hypercube sampling from among the ranges of synaptic weights that maintained biologically realistic firing rates across all populations [29]. In that work, in which no learning occurred, however, the cortico-striatal projections to the action channels were considered to be independent, and variability in decisions was introduced by directly and independently biasing the corticostriatal weights associated with each option, resulting in networks that were inherently biased towards one of the choices.

Hence, some sampled network configurations were biased towards one of the choices. Because we studied the evolution of the control ensembles under plasticity in this work, to ensure that any effects observed during learning were not confounded by pre-existing biases, we started with completely unbiased networks. Hence we resampled the networks from the joint synaptic weight distribution using genetic algorithms (see below) and isolated 300 networks that produced firing rates of all CBGT populations within experimentally observed ranges. The comparison of the ranges for free parameters in [29] and the present work are listed in S2 Table. The actual firing rate distributions are shown in S1 FigA. The networks before plasticity showed a diversity of reaction times (RTs, S1 FigB). The RT distribution was divided into 3 equal tertiles and used to define “fast” (orange), “intermediate” (brown) and “slow” (red) networks. All of the networks before plasticity showed chance levels of accuracy (S1 FigC).

## Dopamine-dependent plasticity of corticostriatal synaptic weights

Synaptic plasticity in CBGTPy is implemented using a dopamine-dependent plasticity rule, where the synaptic updates are governed solely by local factors, without requiring individual neurons to access information about the global system state. This rule is an adaptation of the plasticity mechanism presented in [31,34].

At each corticostriatal AMPA synapse, the model tracks three key values: eligibility  $E(t)$ , weight  $w(t)$ , and conductance  $g_x(t)$ . The conductance is associated with the synaptic current. How much the conductance grows with each pre-synaptic spike is determined by the weight. The weight is the plastic element in the system, which changes over time depending on the time courses of eligibility and dopamine release.

At a computational level,  $E(t)$ , representing a synapse’s eligibility to undergo weight modification, depends on the relative spike times of the pre- and post-synaptic neurons involved in the synapse. To compute this quantity, we first define the variables  $A_{PRE}(t)$  and  $A_{POST}(t)$ , which serve as instantaneous estimates of the recent levels of pre- and post-synaptic spiking, respectively. Each time a spike occurs in the pre- or post-synaptic cell, these values are increased by a fixed amount ( $\Delta_{PRE}$  and  $\Delta_{POST}$ , respectively), and between spikes, they decay exponentially with a time decay constant  $\tau_{PRE}$  and  $\tau_{POST}$ , respectively. That is,

$$\begin{aligned}\frac{dA_{PRE}}{dt} &= \frac{1}{\tau_{PRE}} (\Delta_{PRE} X_{PRE}(t) - A_{PRE}(t)), \\ \frac{dA_{POST}}{dt} &= \frac{1}{\tau_{POST}} (\Delta_{POST} X_{POST}(t) - A_{POST}(t))\end{aligned}$$

where  $X_{PRE}(t)$  and  $X_{POST}(t)$  are sums of Dirac delta functions representing the spike trains of the two neurons. These sums take the form

$$X_{PRE} = \sum_{t_s \in \mathcal{X}_{Cx}} \delta(t - t_s), \quad X_{POST} = \sum_{t_s \in \mathcal{X}_{SPN}} \delta(t - t_s),$$

where  $t_s$  is the spike onset time,  $\mathcal{X}_{Cx}$  is the set of times of all spikes of the cortical neuron involved in the synapse, and  $\mathcal{X}_{SPN}$  refers to the set of spike times of the target postsynaptic neuron within the striatum. Eligibility ( $E(t)$ ) changes over time according to

$$\frac{dE}{dt} = \frac{1}{\tau_E} (X_{POST}(t) A_{PRE}(t) - X_{PRE}(t) A_{POST}(t) - E) \quad (1)$$

where  $\tau_E$  is a time constant. The corticostriatal synaptic conductance  $g_x$  takes the value of the synaptic weight,  $w(t)$ , at each pre-synaptic spike time and decays exponentially in-between these spikes:

$$\frac{dg_x}{dt} = \sum_{t_s \in \mathcal{X}_{Cx}} w(t_s) \delta(t - t_s) - \frac{g_x}{\tau_{AMPA}},$$

where  $x$  labels the specific connection,  $\tau_{AMPA}$  is the decay time constant associated with AMPA synapses, and  $w$  itself changes over time based on dopamine release and the post-synaptic neuron's eligibility. The evolution of  $w$  is given by

$$\frac{dw}{dt} = [\alpha_w^j E(t) f(K_{DA})(w_{max}^j - w)]^+ + [\alpha_w^j E(t) f(K_{DA})(w - w_{min}^j)]^-, \quad (2)$$

where  $[\cdot]^+$  ( $[\cdot]^-$ ) represents a function whose output is the value inside the brackets if it is positive (negative) and 0 otherwise. The learning rate is denoted in Eq (2) by  $\alpha_w^j$ , for  $j \in \{dSPN, iSPN\}$ , depending on to which of the two populations the post-synaptic neuron belongs. This rate has a positive sign for dSPN neurons and a negative one for iSPN neurons to reproduce the observation that positive feedback signals lead to a strengthening of the eligible direct pathway connections and a weakening of the eligible indirect pathway connections. Furthermore,  $w_{max}^j$  and  $w_{min}^j$  are upper and lower bounds for the weight  $w$ , respectively, for  $j \in \{dSPN, iSPN\}$ .

In Eq (2), the variable  $K_{DA}$  represents the level of available dopamine in the network, which is computed from the amount of dopamine released via the differential equation

$$\frac{dK_{DA}}{dt} = C_{scale} \sum_j (DA_{inc}(t_j) - K_{DA}) \delta(t_j) - \frac{K_{DA}}{\tau_{DA}},$$

where  $DA_{inc}(t_j)$  the increment of dopamine, relative to a baseline level, that is delivered at time  $t_j$ . That is, after a specific decision  $i$  is made at time  $t_j$ , a reward value  $r_i(t_j)$  associated to action  $i$  is received, which induces a dopamine increment based on the reward prediction error

$$DA_{inc}(t_j) = r_i(t_j) - Q_i(t_j),$$

where  $Q_i(t_j)$  is the expected reward for action  $i$  at time  $t_j$ . This expected reward obeys the update rule

$$Q_i(t_j^+) = Q_i(t_j) + \alpha_Q (r_i(t_j) - Q_i(t_j)),$$

where  $\alpha_Q \in [0, 1]$  is the learning rate for action values and  $Q(t_j^+) = \lim_{t \rightarrow t_j^+} Q(t)$ . The update of  $DA_{inc}$  in turn impacts the evolution of  $K_{DA}$ . Finally, the function  $f(K_{DA})$  in Eq (2) represents the impact that the available dopamine  $K_{DA}$  has on plasticity, such that, if the target neuron lies in the dSPN population, then

$$f(K_{DA}) = \begin{cases} -\gamma, & \text{if } K_{DA} < -\mu, \\ \frac{\gamma}{\mu} K_{DA}, & \text{if } K_{DA} \geq -\mu, \end{cases}$$

while if the target neuron lies in the iSPN population, then

$$f(K_{DA}) = \begin{cases} \varepsilon \frac{\gamma}{\mu} K_{DA}, & \text{if } K_{DA} < \mu, \\ \varepsilon \gamma, & \text{if } K_{DA} \geq \mu. \end{cases}$$

for fixed, positive scaling parameters  $\gamma, \mu$ . Parameters values used for the plasticity implementation can be found in [S3 Table](#)

## Genetic algorithms

We use genetic algorithms to sample a high-dimensional (14-d) space of parameters and choose a collection of networks that operate in a cortico-basal ganglia-thalamic driven regime and exhibit a heterogeneous range of response times. To this end, we used the DEAP library [78] to run a genetic algorithm (GA) designed to sample CBGT networks with parameters from the ranges used in our previous work [29]. Two additional criteria were used for the optimization function of the GA, namely (a) the network should produce trial timeouts (when no action was selected within 1000 ms) on fewer than 1% of trials, and (b) the network should be cortico-basal-ganglia driven; that is, the correlation between cortical activity and striatal activity should be positive. The first criterion ensured that we had ample decision trials included in the data, as needed to appropriately fit the DDM parameters (timeouts are dropped before fitting the DDM parameters). The second criterion ensured that the networks did not operate in a cortico-thalamic driven regime, in which cortical inputs alone directly pushed thalamic firing over the decision threshold.

The range for each parameter specified in past work [29] was divided into 30 bins and this grid was sampled to create populations. The indices of each bin served as a pointer to the actual values of the parameters in the ranges considered. The GA uses these indices to create, mate and mutate the populations. This ensures that the values of parameters remain within their specified ranges. For example, suppose that parameter  $A$  has range  $(-2.0, 2.0)$  and parameter  $B$  has range  $(-0.3, 1.0)$  and these ranges are each divided into 5 bins. The grids for parameters  $A$  and  $B$  would then be:

$$A_{grid} = (-2 \quad -1 \quad 0 \quad 1 \quad 2)$$

$$B_{grid} = (-0.3 \quad 0.025 \quad 0.35 \quad 0.675 \quad 1).$$

If individual population members have indices  $ind_1 = (0 \ 1)$  and  $ind_2 = (4 \ 0)$  for  $(A, B)$ , then they have  $(A, B) = (-2, 0.025)$  and  $(A, B) = (2, -0.3)$ , respectively. Supposed that the individuals mate by crossing over the 1st and 2nd elements. Then  $ind_3 = (4 \ 1)$  with parameter values  $(2, 0.025)$  and  $ind_4 = (0 \ 0)$  with parameter values  $(-2, -0.3)$ . The individuals  $ind_3$  and  $ind_4$  are included in the next iteration of evolution.

New individuals created from mating were used to overwrite the original individuals that were mated together (*cxSimulatedBinary*). The individuals could also mutate by shuffling of the indices of the attributes (*mutShuffleIndexes*) with a probability of 0.2. After a round of mating and mutation, tuples of two values for each individual, namely the % of timeouts and the Pearson's correlation coefficient between cortical and striatal activity, were compared to select the individuals for the next round of evolution. The selection algorithm that was used was tournament selection (*selTournament*) of size 3, which picked the best individual among 3 randomly chosen individuals, 10 times, where 10 is the size of the population of networks in every iteration of the GA. During every iteration, any network configuration that met the criteria (a) and (b) above was saved as a correct solution. The GA was run for 2000 iterations or until 300 solutions were found, whichever was sooner. Post hoc, we confirmed that the firing rates of the members of the final, selected populations remained within the originally targeted ranges ([S1 Fig](#)).

## DDM fits

A drift diffusion model (DDM) as shown in [Fig 3](#) assumes decision to be a noisy random walk with a drift rate ( $v$ ) towards one of the boundaries that represent the two choices. The difference between the two boundary heights ( $a$ ) represents the decision threshold required to commit to a decision. The amount of time required for the decision variable to reach one of the boundaries constitute the reaction times (RT). The DDM also has additional parameters, non decision time ( $t$ ), that may correspond to the sensory processing time or other processes not related to the decision time and the inherent bias

(z) towards one of the choices.

$$dx(t) = v dt + \sigma dW(t)$$

where  $x(t)$  represents the noisy decision variable,  $v$  represents drift rate,  $\sigma$  represents the level of noise in the process,  $W(t)$  represents a *Weiner* process and was drawn from a normal distribution with mean zero and variance  $dt$ . The process  $x(t)$  terminated at the minimum time  $t$ , when  $x(t)$  reaches one of the boundaries (i.e.  $x(t) = 0$  or  $x(t) = a$ ).

The DDM parameters were fit to each of the 300 selected networks independently using the HDDM package [35]. In order to ensure that the HDDM fits describe the choice and RT distributions well, we compared the post-predictive distributions of the network simulations with those generated by the corresponding DDM parameters, both before (S2 FigA) and after (S2 FigC) plasticity. The quantile-quantile plots for percentiles 5, 10, ..., 90, 95 show a significant and very high correlation between the network-generated and DDM-generated data (S2 FigB and S2 FigD).

### Accuracy, RT and reward rate manifolds

The manifolds shown in Fig 2A–2C were generated by simulating the DDM with all combinations of drift rate ( $v$ ) and boundary height ( $a$ ) values that the naive CBGT networks can show before and after plasticity. The values of RTs, accuracy and reward rates were averaged over 15 seeds of 200 trials each.

### Upward mapping

The DDM parameters and activity of the CBGT nuclei for our 300 network configurations, before plasticity, were used to identify CBGT control ensembles through canonical correlation analysis (CCA), as was also done in our previous work [29] and is illustrated in Fig 3. The CCA identifies independent components in two sets of variables such that the correlation between the projections of the data onto these components is maximized. It is used to uncover linear combinations of features that have the strongest correlations with each other. In our previous work [29], we applied CCA to link behavior (captured by DDM parameters) and CBGT network activity (described by the sums and differences of activity levels in corresponding populations across channels). This analysis identified “control ensembles”: linear combinations of network features that modulate decisions in a two-choice task. Both our previous and current studies consistently reveal three key control ensembles, corresponding to linear combinations that we term responsiveness, pliancy, and choice. These are described in more detail in the section *Low-dimensional control ensembles that map to general decision policies*.

The CCA scores were calculated using  $k$ -fold validation ( $k=4$ ), where the 300 networks were divided into groups of 4 (75 networks each) and a CCA score was calculated for each of the groups. The CCA scores for actual data were compared with a shuffled version of data (firing rates and DDM components for 300 networks) and the set of components giving rise to the maximum CCA score, which we found to include three elements as in our previous work [29], were selected.

### Modulation of control ensembles by plasticity

We used a single approach to compute a set of effective drivers of the control ensembles either from the full collection of CBGT networks or from one of the network subtypes (fast, intermediate, or slow) that we considered. Let  $X \in \{\text{all}, \text{fast}, \text{intermediate}, \text{slow}\}$  denote the class of networks being used. From the set of vectors of changes in CBGT firing rates computed by subtracting firing rates before plasticity from those after plasticity ( $\Delta F_X$ ), we extracted 5 principal components (PCs) that together explain at least about 90% of the variance (Fig 4A and S7 FigA).  $\Delta F_X$  was then projected onto these 5 PCs to form the target matrix  $P_X$ . Specifically, we computed

$$P_X = (\Delta F_X) V_X \quad (3)$$

where the 5 PCs comprise the columns of  $V_X$ . Note that  $P_X$  is an  $n$  by 5 matrix, where  $n$  is the number of firing rate data vectors used.  $\Delta F_X$  was also projected onto the three control ensemble components obtained from the full collection of baseline networks before plasticity, via the mapping

$$C_X = (\Delta F_X)U \quad (4)$$

where the components of the 3 control ensembles form the columns of  $U$ , such that  $C_X$  is an  $n$  by 3 matrix. Finally, we found the least squares solution  $S_X$ , representing the element in the range of  $C_X$  that is closest to  $P_X$ , from the normal equation

$$S_X = (C_X^T C_X)^{-1} C_X^T P_X. \quad (5)$$

The least squares solution  $S_X$  is a  $3 \times 5$  matrix independent of  $n$ . The columns of  $S_{all}$  are displayed in Fig 4B. The sums of the columns of the appropriate  $S_X$ , each weighted by the percent of variance explained, comprise Figs 5C and S10 Fig ( $X = \text{fast}$ ,  $X = \text{intermediate}$ , and  $X = \text{slow}$ ), as well as Figs 5A and S7 FigB ( $X = \text{all}$ ).

## Reward rates

The reward rate was calculated as:

$$\begin{aligned} RR &= \frac{1 - p(\text{err})}{DT + T_0} \\ &= \frac{\text{accuracy}}{RT} \end{aligned}$$

where  $p(\text{err})$  denotes the error rate and where the reaction time,  $RT$ , is the sum of the decision time,  $DT$ , and the additional non-decision time that arises within each trial,  $T_0$ , which in our analysis is ascribed to the onset delay represented by the DDM parameter  $t$ .

## Plasticity stages

The effect of plasticity on the network was studied at four stages: (a) after 2 trials of plasticity, (b) after 2 additional trials (total 4) of plasticity, (c) after 2 more additional trials (total 6) of plasticity, (d) after 9 additional trials (total 15) of plasticity. The state of the network was frozen at each of these stages by suspending the plasticity, so that we could use the frozen network to perform probe trials. The choices and reaction times from the probe trials were used to calculate DDM parameters and reward rate distributions for each stage of plasticity, based on upward mapping and CCA, and thus to generate the trajectories in Fig 2, the time courses in S9 Fig, and the 2-trial results in Figs 5, S10 Fig, and S11 Fig.

## Supporting information

**S1 Table. Percentage of first pairs of trials for which networks encounter each possible reward sequence.** Slow networks encounter a higher proportion of two consecutively unrewarded choices (U-U) and fewer R-U sequences than intermediate and fast networks.

(TEX)

**S2 Table. Comparison of parameter ranges used for Latin hypercube sampling (LHS) in [29] and genetic algorithms (GA) in the present work.**

(TEX)

**S3 Table. Plasticity model parameters.**

(TEX)

**S4 Table. External current parameters.** Parameters used to describe the external current ( $I_{ext}$ ) arriving at the different populations of the CBGT network. From the third column to the last, we specify the different parameters used to describe the external current impinging in each population specified in column 1 and for the specific type of receptors. A non described receptor type means that the parameters are considered to be zero. The time decay constant  $\tau$  is the same for all populations and only depends on the type of receptor being  $\tau = 2\text{ ms}$  if the receptor type is AMPA and  $\tau = 5\text{ ms}$  if it is GABA. \* Values in this row are the ones used when no intrinsic separation of neurons is considered.

(TEX)

**S5 Table. Firing rates observed in CBGT nuclei.** The second column refers to the firing frequency ranges observed experimentally during baseline for each population set in the first column, whereas the third column refers to the ranges observed during decision tasks. In both cases, the ranges reflect experimental data from primates and rats (see references in the last column).

(TEX)

**S1 Fig. Network firing rates, RTs, and accuracy before plasticity.** A: The distributions of average firing rates for the 9 CBGT regions based on 300 networks. An average was calculated for each population over the whole simulation time. One example each from three categories of network – fast (orange), intermediate (brown) and slow (red) – are marked on the distribution. B: The networks before plasticity were categorized as fast, intermediate and slow based on a tertile split of the reaction time (RT) distribution (vertical dashed linebs). The RTs for the exemplar fast (orange), intermediate (brown) and slow (red) networks are marked. C: The average accuracies of all 300 networks. The accuracy distribution is centered around 50% (0.5) because the networks had not yet undergone plasticity.

(TIFF)

**S2 Fig. DDM fits for 300 networks before and after plasticity.** A: The post-predictive choice (i.e., split between positive and negative RTs) and RT distributions from the naive (before plasticity) network simulations (red, “actual”) and distributions generated by the DDM parameters fitted to the data (purple, “generated”) separately for fast, intermediate and slow networks. Note the near-symmetry of the two RT peaks for the two choices (left → positive, right → negative). B: Quantile-Quantile plots for the distributions shown in A for percentiles in steps of 5 (i.e., 5, 10...90, 95). The Pearson correlation and p-value between the actual and generated data are annotated in green. The Pearson correlation was significant for all three network types (0.94, 0.92 and 1.0 for fast, intermediate and slow networks, respectively). C: Same as A but after plasticity with the left choice (positive RTs) rewarded. D: Same as B but after plasticity.

(TIFF)

**S3 Fig. Comparison of DDM and behavioral measures for all 300 networks before (blue) and after (pink) plasticity.** The subplots on the diagonal represent the marginal distributions for DDM parameters ( $a$ ,  $t$ ,  $v$ ) and behavioral features (RT and accuracy). The onset delay ( $t$ ) shows a decrease, the drift rate ( $v$ ) shows an increase, RTs show a decrease, and accuracy shows an increase after plasticity. The off-diagonal subplots show the pairwise covariances.

(TIFF)

**S4 Fig. Evolution of behavioral measures for 300 networks over 16 trials with plasticity.** A: Network behavior was assessed after each of 2, 4, 6, 9 and 15 trials. The RTs steadily decreased for all three network categories: fast (orange), intermediate (brown) and slow (red). The average over all 300 networks also showed a steady decrease as shown in black markers and lines. B: The accuracy for the three categories of the networks and the average over all 300 networks increased with plasticity. C: The reward rate for three categories of network and the average over 300 networks increased with plasticity. D: The distribution of differences in cosine distance, measured relative to the direction of greatest increase,

for changes in RT vs accuracy, RT vs reward rate, and accuracy vs reward rate for all 300 networks and all stages of plasticity. The comparisons with reward rate yield distributions skewed to significantly above 0, suggesting that the cosine distances are lowest for reward rates. E: Absolute cosine distance distributions shown separately for the three network classes, fast (orange), intermediate (brown) and slow (red).

(TIFF)

**S5 Fig. Distance from the optimal direction for different levels of reward probabilities.** A) DDM ( $a$ ,  $v$ ) and behavioral (RT, accuracy) parameters for reward probabilities: 100% (red), 95% (pink), 90% (yellow) and 75% (khaki) measured for 150 networks. The distributions before plasticity are shown in blue. As the conflict increases, accuracy decreases and RTs show a lower average decrease after plasticity B) Cosine distances with respect to the RT (minimization), Accuracy (maximization) and Reward rate (maximization) vectors for the three reward probabilities. The cosine distances scale with the reward probabilities but remain the lowest for reward rates, for 100% reward probabilities and 150 networks in total:  $F[449,2] = 34.8$ ,  $p < 0.00001$ , for 95%:  $F[449,2] = 89.9$ ,  $p < 0.00001$ , for 90%:  $F[449,2] = 88.8$ ,  $p < 0.00001$  and for 75%:  $F[449,2] = 25.03$ ,  $p < 0.00001$ . C) The distribution of differences in cosine distance, measured relative to the direction of greatest increase, for changes in RT vs accuracy (100%:  $t(149) = 3.35$ ,  $p = 0.001$ ; 95%:  $t(149) = 6.72$ ,  $p < 0.0001$ ; 90%:  $t(149) = 8.9$ ,  $p < 0.0001$ ; 75%:  $t(149) = 7.93$ ,  $p < 0.0001$ ), RT vs reward rate (100%:  $t(149) = 9.65$ ,  $p < 0.0001$ ; 95%:  $t(149) = 12.42$ ,  $p < 0.0001$ ; 90%:  $t(149) = 13.17$ ,  $p < 0.0001$ ; 75%:  $t(149) = 9.9$ ,  $p < 0.0001$ ), and accuracy vs reward rate (100%:  $t(149) = 6.07$ ,  $p < 0.0001$ ; 95%:  $t(149) = 8.36$ ,  $p < 0.0001$ ; 90%:  $t(149) = 4.06$ ,  $p < 0.0001$ ; 75%:  $t(149) = 1.02$ ,  $p = 0.3$ ) for all the reward probabilities. The comparisons with reward rate yield distributions skewed to significantly above 0 (except for 75%), suggesting that the cosine distances are lowest for reward rates for all the reward probabilities.

(TIFF)

**S6 Fig. Comparison with RL-DDM.** Deviance Information Criterion (DIC) values are compared for DDM and RL-DDM fits across three phases: before, during, and after plasticity. Across all network types (Fast, Intermediate, and Slow) the DDM consistently outperforms the RL-DDM before and after plasticity. However, during plasticity, both models yield comparable fits, indicating that RL-DDM is comparable to DDM in capturing network behavior when plasticity is ongoing, but not once it stabilizes.

(TIFF)

**S7 Fig. The least squares solution  $S$  pooled over the network types.** A:) Cumulative variance explained by the first 10 principal components (PC) derived from the changes in firing rates from before to after plasticity. The dashed line indicates 90% of variance explained. The analysis was done for all the networks pooled together (blue line) and separately for fast (orange), intermediate (brown) and slow (red) networks. For all networks pooled together as well as the separated slow and intermediate networks, the first 5 PCs explain more than 90% of the variance, whereas for fast networks 1 PC suffices. B: The weighted sum of the columns of  $S$  (see main text, Fig 4B), pooled over all three network classes (fast, intermediate and slow), shows that the observed changes in firing rates correspond to increased loadings of the responsiveness, pliancy and choice ensembles of the CBGT network, to differing extents.

(TIFF)

**S8 Fig. Reconstruction of firing rate changes from the least squares solution  $S$  for the three network classes.**

(A) The first 5 PCs for the firing rate changes in the fast networks. Although the 1st PC explains around 90% of the variance for fast networks, we used 5 PCs to calculate  $S$  coefficients (Fig 4C) to be consistent with slow and intermediate networks (Supp. S7 FigA). (B,C): Same as (A) for intermediate and slow networks, respectively. (D-F) The dot products of the CCA component vector ( $C$ ) with each of the 5 columns of  $S$ , the least squares solution of  $P = CS$ , provide an approximate reconstruction of the 5 PCs of the changes in firing rate from before to after plasticity, ( $\Delta F$ ). The quality of the reconstruction was checked by projecting  $\Delta F$  onto the original PCs for each network (marked as *Actual* on y-axis) and comparing the results with the projections of  $\Delta F$  onto the reconstructed PCs (marked as *Predicted* on x-axis). The goodness of

fit is calculated as the Spearman rank correlation ( $\rho$ ) between the actual and predicted values. For fast networks (D), the rank correlations ( $\rho$ ) are high and significant ( $p < 0.0001$ ) for all of the PCs as shown, suggesting that the reconstruction is excellent. For intermediate networks (E), the rank correlations are significant for all PCs except the 5th PC. For slow networks (E), the rank correlations are significant for all except 4th and 5th PCs.

(TIFF)

**S9 Fig. Evolution of DDM parameters with plasticity.** (A) The change in boundary height ( $a$ ) due to plasticity is dependent on network type: slow networks (red) show a decrease, intermediate (brown) show little change, and fast (orange) networks show a slight increase. The mean over all networks is shown by large black circles. (B) All network types show a decrease in decision onset time ( $t$ ) due to plasticity. (C) All network types show a strong increase in drift rate ( $v$ ) due to plasticity. (D) Fast networks make more mistakes on average. The histograms show the proportion of unrewarded (“U”) trials encountered by all the three network classes after the first two plasticity trials.

(TIFF)

**S10 Fig. Effect of reward sequences on the weighting coefficients  $S$  for the three network classes.** The weighting coefficients  $S$  shown in Fig 5A combine the three network types. The separated coefficients here show the same trends as the combined ones.

(TIFF)

**S11 Fig. Effect of reward sequences on changes in accuracy and reaction times (RTs).** (A) The change in accuracy showed an increase in all cases, but to different extents. The highest increase in accuracy was for one rewarded and one unrewarded trial (U-R and R-U), due to strengthening of the cortico-striatal projection to dSPNs of the optimal choice along with strengthening of cortico-striatal projections to iSPNs of the sub-optimal choice. (B) The change in RTs after plasticity for the four outcome sequences. All sequences involving at least one rewarded trial yielded a decrease in RT, whereas the sequence with two consecutive unrewarded trials (U-U) induced an increase in RT.

(TIFF)

## Acknowledgments

We thank all members of the exploratory intelligence group for their helpful comments on the manuscript.

## Author contributions

**Conceptualization:** Timothy Verstynen, Jonathan E. Rubin.

**Data curation:** Jyotika Bahuguna.

**Formal analysis:** Jyotika Bahuguna.

**Funding acquisition:** Timothy Verstynen, Jonathan E. Rubin.

**Investigation:** Jyotika Bahuguna.

**Methodology:** Jyotika Bahuguna, Timothy Verstynen, Jonathan E. Rubin.

**Project administration:** Timothy Verstynen, Jonathan E. Rubin.

**Resources:** Timothy Verstynen, Jonathan E. Rubin.

**Software:** Jyotika Bahuguna.

**Supervision:** Timothy Verstynen, Jonathan E. Rubin.

**Validation:** Jyotika Bahuguna, Timothy Verstynen, Jonathan E. Rubin.

**Visualization:** Jyotika Bahuguna, Timothy Verstynen, Jonathan E. Rubin.

**Writing – original draft:** Jyotika Bahuguna, Timothy Verstynen, Jonathan E. Rubin.

**Writing – review & editing:** Jyotika Bahuguna, Timothy Verstynen, Jonathan E. Rubin.

## References

- Gold JI, Shadlen MN. The neural basis of decision making. *Annu Rev Neurosci*. 2007;30:535–74. <https://doi.org/10.1146/annurev.neuro.29.051605.113038> PMID: 17600525
- Cohen JD, McClure SM, Yu AJ. Should I stay or should I go? How the human brain manages the trade-off between exploitation and exploration. *Philos Trans R Soc Lond B Biol Sci*. 2007;362(1481):933–42. <https://doi.org/10.1098/rstb.2007.2098> PMID: 17395573
- Mehlhorn K, Newell BR, Todd PM, Lee MD, Morgan K, Braithwaite VA, et al. Unpacking the exploration–exploitation tradeoff: a synthesis of human and animal literatures. *Decision*. 2015;2(3):191–215. <https://doi.org/10.1037/dec0000033>
- Wilson RC, Bonawitz E, Costa VD, Ebitz RB. Balancing exploration and exploitation with information and randomization. *Curr Opin Behav Sci*. 2021;38:49–56. <https://doi.org/10.1016/j.cobeha.2020.10.001> PMID: 33184605
- Dudman JT, Krakauer JW. The basal ganglia: from motor commands to the control of vigor. *Curr Opin Neurobiol*. 2016;37:158–66. <https://doi.org/10.1016/j.conb.2016.02.005> PMID: 27012960
- Bond K, Dunovan K, Porter A, Rubin JE, Verstynen T. Dynamic decision policy reconfiguration under outcome uncertainty. *Elife*. 2021;10:e65540. <https://doi.org/10.7554/eLife.65540> PMID: 34951589
- Bond K, Rasero J, Madan R, Bahuguna J, Rubin J, Verstynen T. Competing neural representations of choice shape evidence accumulation in humans. *Elife*. 2023;12:e85223. <https://doi.org/10.7554/eLife.85223> PMID: 37818943
- Ratcliff R. A theory of memory retrieval. *Psychological Review*. 1978;85(2):59–108. <https://doi.org/10.1037/0033-295x.85.2.59>
- Ratcliff R, McKoon G. The diffusion decision model: theory and data for two-choice decision tasks. *Neural Comput*. 2008;20(4):873–922. <https://doi.org/10.1162/neco.2008.12-06-420> PMID: 18085991
- Ratcliff R, Smith PL, Brown SD, McKoon G. Diffusion decision model: current issues and history. *Trends Cogn Sci*. 2016;20(4):260–81. <https://doi.org/10.1016/j.tics.2016.01.007> PMID: 26952739
- Bogacz R, Wagenmakers E-J, Forstmann BU, Nieuwenhuis S. The neural basis of the speed-accuracy tradeoff. *Trends Neurosci*. 2010;33(1):10–6. <https://doi.org/10.1016/j.tins.2009.09.002> PMID: 19819033
- Smith PL, Ratcliff R. Psychology and neurobiology of simple decisions. *Trends Neurosci*. 2004;27(3):161–8. <https://doi.org/10.1016/j.tins.2004.01.006> PMID: 15036882
- Bogacz R, Brown E, Moehlis J, Holmes P, Cohen JD. The physics of optimal decision making: a formal analysis of models of performance in two-alternative forced-choice tasks. *Psychol Rev*. 2006;113(4):700–65. <https://doi.org/10.1037/0033-295X.113.4.700> PMID: 17014301
- Leng X, Yee D, Ritz H, Shenhav A. Dissociable influences of reward and punishment on adaptive cognitive control. *PLoS Comput Biol*. 2021;17(12):e1009737. <https://doi.org/10.1371/journal.pcbi.1009737> PMID: 34962931
- Zacksenhouse M, Bogacz R, Holmes P. Robust versus optimal strategies for two-alternative forced choice tasks. *J Math Psychol*. 2010;54(2):230–46. <https://doi.org/10.1016/j.jmp.2009.12.004> PMID: 23180885
- Yttri EA, Dudman JT. Opponent and bidirectional control of movement velocity in the basal ganglia. *Nature*. 2016;533(7603):402–6. <https://doi.org/10.1038/nature17639> PMID: 27135927
- Tecuapetla F, Jin X, Lima SQ, Costa RM. Complementary contributions of striatal projection pathways to action initiation and execution. *Cell*. 2016;166(3):703–15. <https://doi.org/10.1016/j.cell.2016.06.032> PMID: 27453468
- Herz DM, Bange M, Gonzalez-Escamilla G, Auer M, Ashkan K, Fischer P, et al. Dynamic control of decision and movement speed in the human basal ganglia. *Nat Commun*. 2022;13(1):7530. <https://doi.org/10.1038/s41467-022-35121-8> PMID: 36476581
- de A Marcelino AL, Gray O, Al-Fatly B, Gilmour W, Douglas Steele J, Kühn AA, et al. Pallidal neuromodulation of the explore/exploit trade-off in decision-making. *Elife*. 2023;12:e79642. <https://doi.org/10.7554/eLife.79642> PMID: 36727860
- Geddes CE, Li H, Jin X. Optogenetic editing reveals the hierarchical organization of learned action sequences. *Cell*. 2018;174(1):32–43.e15. <https://doi.org/10.1016/j.cell.2018.06.012> PMID: 29958111
- Albin RL, Young AB, Penney JB. The functional anatomy of disorders of the basal ganglia. *Trends in Neurosciences*. 1995;18(2):63–4. [https://doi.org/10.1016/0166-2236\(95\)80020-3](https://doi.org/10.1016/0166-2236(95)80020-3)
- DeLong MR. Primate models of movement disorders of basal ganglia origin. *Trends Neurosci*. 1990;13(7):281–5. [https://doi.org/10.1016/0166-2236\(90\)90110-v](https://doi.org/10.1016/0166-2236(90)90110-v) PMID: 1695404
- Lo C-C, Wang X-J. Cortico-basal ganglia circuit mechanism for a decision threshold in reaction time tasks. *Nat Neurosci*. 2006;9(7):956–63. <https://doi.org/10.1038/nn1722> PMID: 16767089

24. Bogacz R, Gurney K. The basal ganglia and cortex implement optimal decision making between alternative actions. *Neural Comput.* 2007;19(2):442–77. <https://doi.org/10.1162/neco.2007.19.2.442> PMID: 17206871
25. Bogacz R, Martin Moraud E, Abdi A, Magill PJ, Baufreton J. Properties of neurons in external globus pallidus can support optimal action selection. *PLoS Comput Biol.* 2016;12(7):e1005004. <https://doi.org/10.1371/journal.pcbi.1005004> PMID: 27389780
26. Dunovan K, Verstynen T. Believer-skeptic meets actor-critic: rethinking the role of basal ganglia pathways during decision-making and reinforcement learning. *Front Neurosci.* 2016;10:106. <https://doi.org/10.3389/fnins.2016.00106> PMID: 27047328
27. Bariselli S, Fobbs WC, Creed MC, Kravitz AV. A competitive model for striatal action selection. *Brain Res.* 2019;1713:70–9. <https://doi.org/10.1016/j.brainres.2018.10.009> PMID: 30300636
28. Dunovan K, Vich C, Clapp M, Verstynen T, Rubin J. Reward-driven changes in striatal pathway competition shape evidence evaluation in decision-making. *PLoS Comput Biol.* 2019;15(5):e1006998. <https://doi.org/10.1371/journal.pcbi.1006998> PMID: 31060045
29. Vich C, Clapp M, Rubin JE, Verstynen T. Identifying control ensembles for information processing within the cortico-basal ganglia-thalamic circuit. *PLoS Comput Biol.* 2022;18(6):e1010255. <https://doi.org/10.1371/journal.pcbi.1010255> PMID: 35737720
30. Frank MJ. Linking across levels of computation in model-based cognitive neuroscience. An introduction to model-based cognitive neuroscience. 2015. p. 159–77.
31. Vich C, Dunovan K, Verstynen T, Rubin J. Corticostriatal synaptic weight evolution in a two-alternative forced choice task: a computational study. *Communications in Nonlinear Science and Numerical Simulation.* 2020;82:105048. <https://doi.org/10.1016/j.cnsns.2019.105048>
32. Masís J, Chapman T, Rhee JY, Cox DD, Saxe AM. Strategically managing learning during perceptual decision making. *Elife.* 2023;12:e64978. <https://doi.org/10.7554/eLife.64978> PMID: 36786427
33. Robbins H. Some aspects of the sequential design of experiments. *Bull Amer Math Soc.* 1952;58(5):527–35. <https://doi.org/10.1090/s0002-9904-1952-09620-8>
34. Clapp M, Bahuguna J, Giossi C, Rubin JE, Verstynen T, Vich C. CBGTPy: an extensible cortico-basal ganglia-thalamic framework for modeling biological decision making. *bioRxiv.* 2024;:2023.09.05.556301. <https://doi.org/10.1101/2023.09.05.556301> PMID: 37732280
35. Wiecki TV, Sofer I, Frank MJ. HDDM: Hierarchical Bayesian estimation of the drift-diffusion model in Python. *Frontiers in Neuroinformatics.* 2013;:14.
36. Fengler A, Bera K, Pedersen ML, Frank MJ. Beyond drift diffusion models: fitting a broad class of decision and reinforcement learning models with HDDM. *J Cogn Neurosci.* 2022;34(10):1780–805. [https://doi.org/10.1162/jocn\\_a\\_01902](https://doi.org/10.1162/jocn_a_01902) PMID: 35939629
37. Bowman NE, Kording KP, Gottfried JA. Temporal integration of olfactory perceptual evidence in human orbitofrontal cortex. *Neuron.* 2012;75(5):916–27. <https://doi.org/10.1016/j.neuron.2012.06.035> PMID: 22958830
38. Shadlen MN, Kiani R. Decision making as a window on cognition. *Neuron.* 2013;80(3):791–806. <https://doi.org/10.1016/j.neuron.2013.10.047> PMID: 24183028
39. Malhotra G, Leslie DS, Ludwig CJH, Bogacz R. Time-varying decision boundaries: insights from optimality analysis. *Psychon Bull Rev.* 2018;25(3):971–96. <https://doi.org/10.3758/s13423-017-1340-6> PMID: 28730465
40. Ratcliff R, Frank MJ. Reinforcement-based decision making in corticostriatal circuits: mutual constraints by neurocomputational and diffusion models. *Neural Comput.* 2012;24(5):1186–229. [https://doi.org/10.1162/NECO\\_a\\_00270](https://doi.org/10.1162/NECO_a_00270) PMID: 22295983
41. Rubin J, Lee DD, Sompolinsky H. Equilibrium properties of temporally asymmetric Hebbian plasticity. *Phys Rev Lett.* 2001;86(2):364–7. <https://doi.org/10.1103/PhysRevLett.86.364> PMID: 11177832
42. Güti R, Aharonov R, Rotter S, Sompolinsky H. Learning input correlations through nonlinear temporally asymmetric Hebbian plasticity. *J Neurosci.* 2003;23(9):3697–714. <https://doi.org/10.1523/JNEUROSCI.23-09-03697.2003> PMID: 12736341
43. Williams RJ. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning.* 1992;8(3–4):229–56. <https://doi.org/10.1023/a:1022672621406>
44. Mahadevan S. Average reward reinforcement learning: foundations, algorithms, and empirical results. *Machine Learning.* 1996;22(1–3):159–95. <https://doi.org/10.1023/a:1018064306595>
45. Sutton RS, Barto AG. Reinforcement learning: an introduction. Cambridge: MIT Press; 1998.
46. Loewenstein Y, Seung HS. Operant matching is a generic outcome of synaptic plasticity based on the covariance between reward and neural activity. *Proc Natl Acad Sci U S A.* 2006;103(41):15224–9. <https://doi.org/10.1073/pnas.0505220103> PMID: 17008410
47. Hikosaka O, Rand MK, Miyachi S, Miyashita K. Learning of sequential movements in the monkey: process of learning and retention of memory. *J Neurophysiol.* 1995;74(4):1652–61. <https://doi.org/10.1152/jn.1995.74.4.1652> PMID: 8989401
48. Balci F, Simen P, Niyogi R, Saxe A, Hughes JA, Holmes P, et al. Acquisition of decision making criteria: reward rate ultimately beats accuracy. *Atten Percept Psychophys.* 2011;73(2):640–57. <https://doi.org/10.3758/s13414-010-0049-7> PMID: 21264716
49. Dutilh G, Vandekerckhove J, Tuerlinckx F, Wagenmakers E-J. A diffusion model decomposition of the practice effect. *Psychon Bull Rev.* 2009;16(6):1026–36. <https://doi.org/10.3758/16.6.1026> PMID: 19966251
50. Grahek I, Leng X, Musslick S, Shenhav A. Control adjustment costs limit goal flexibility: Empirical evidence and a computational account. *bioRxiv.* 2025;:2023.08.22.554296. <https://doi.org/10.1101/2023.08.22.554296> PMID: 37662382
51. Don HJ, Otto AR, Cornwall AC, Davis T, Worthy DA. Learning reward frequency over reward probability: A tale of two learning rules. *Cognition.* 2019;193:104042. <https://doi.org/10.1016/j.cognition.2019.104042> PMID: 31430606

52. Frank MJ, Santamaria A, O'Reilly RC, Willcutt E. Testing computational models of dopamine and noradrenaline dysfunction in attention deficit/hyperactivity disorder. *Neuropsychopharmacology*. 2007;32(7):1583–99. <https://doi.org/10.1038/sj.npp.1301278> PMID: 17164816
53. Hamid AA, Pettibone JR, Mabrouk OS, Hetrick VL, Schmidt R, Vander Weele CM, et al. Mesolimbic dopamine signals the value of work. *Nat Neurosci*. 2016;19(1):117–26. <https://doi.org/10.1038/nn.4173> PMID: 26595651
54. Fontanesi L, Gluth S, Spektor MS, Rieskamp J. A reinforcement learning diffusion decision model for value-based decisions. *Psychon Bull Rev*. 2019;26(4):1099–121. <https://doi.org/10.3758/s13423-018-1554-2> PMID: 30924057
55. Pedersen ML, Frank MJ, Biele G. The drift diffusion model as the choice rule in reinforcement learning. *Psychon Bull Rev*. 2017;24(4):1234–51. <https://doi.org/10.3758/s13423-016-1199-y> PMID: 27966103
56. Miletic S, Boag RJ, Trutti AC, Stevenson N, Forstmann BU, Heathcote A. A new model of decision processing in instrumental learning tasks. *Elife*. 2021;10:e63055. <https://doi.org/10.7554/eLife.63055> PMID: 33501916
57. Uehara S, Mawase F, Therrien AS, Cherry-Allen KM, Celnik P. Interactions between motor exploration and reinforcement learning. *J Neurophysiol*. 2019;122(2):797–808. <https://doi.org/10.1152/jn.00390.2018> PMID: 31242063
58. Liquin EG, Gopnik A. Children are more exploratory and learn more than adults in an approach-avoid task. *Cognition*. 2022;218:104940. <https://doi.org/10.1016/j.cognition.2021.104940> PMID: 34715584
59. Hanks T, Kiani R, Shadlen MN. A neural mechanism of speed-accuracy tradeoff in macaque area LIP. *Elife*. 2014;3:e02260. <https://doi.org/10.7554/eLife.02260> PMID: 24867216
60. Banaie Boroujeni K, Oemisch M, Hassani SA, Womelsdorf T. Fast spiking interneuron activity in primate striatum tracks learning of attention cues. *Proc Natl Acad Sci U S A*. 2020;117(30):18049–58. <https://doi.org/10.1073/pnas.2001348117> PMID: 32661170
61. Pasupathy A, Miller EK. Different time courses of learning-related activity in the prefrontal cortex and striatum. *Nature*. 2005;433(7028):873–6. <https://doi.org/10.1038/nature03287> PMID: 15729344
62. Cavanagh JF, Wiecki TV, Cohen MX, Figueroa CM, Samanta J, Sherman SJ, et al. Subthalamic nucleus stimulation reverses mediofrontal influence over decision threshold. *Nat Neurosci*. 2011;14(11):1462–7. <https://doi.org/10.1038/nn.2925> PMID: 21946325
63. Zaghoul KA, Weidemann CT, Lega BC, Jaggi JL, Baltuch GH, Kahana MJ. Neuronal activity in the human subthalamic nucleus encodes decision conflict during action selection. *J Neurosci*. 2012;32(7):2453–60. <https://doi.org/10.1523/JNEUROSCI.5815-11.2012> PMID: 22396419
64. Frank MJ, Scheres A, Sherman SJ. Understanding decision-making deficits in neurological conditions: insights from models of natural action selection. *Philos Trans R Soc Lond B Biol Sci*. 2007;362(1485):1641–54. <https://doi.org/10.1098/rstb.2007.2058> PMID: 17428775
65. Dahlin E, Bäckman L, Neely AS, Nyberg L. Training of the executive component of working memory: subcortical areas mediate transfer effects. *Restor Neurol Neurosci*. 2009;27(5):405–19. <https://doi.org/10.3233/RNN-2009-0492> PMID: 19847067
66. Tremblay L, Hollerman JR, Schultz W. Modifications of reward expectation-related neuronal activity during learning in primate striatum. *J Neurophysiol*. 1998;80(2):964–77. <https://doi.org/10.1152/jn.1998.80.2.964> PMID: 9705482
67. Murayama K, Matsumoto M, Izuma K, Matsumoto K. Neural basis of the undermining effect of monetary reward on intrinsic motivation. *Proc Natl Acad Sci U S A*. 2010;107(49):20911–6. <https://doi.org/10.1073/pnas.1013305107> PMID: 21078974
68. Shohamy D. Learning and motivation in the human striatum. *Curr Opin Neurobiol*. 2011;21(3):408–14. <https://doi.org/10.1016/j.conb.2011.05.009> PMID: 21658933
69. Peters S, Crone EA. Increased striatal activity in adolescence benefits learning. *Nat Commun*. 2017;8(1):1983. <https://doi.org/10.1038/s41467-017-02174-z> PMID: 29259158
70. Yartsev MM, Hanks TD, Yoon AM, Brody CD. Causal contribution and dynamical encoding in the striatum during evidence accumulation. *Elife*. 2018;7:e34929. <https://doi.org/10.7554/eLife.34929> PMID: 30141773
71. Wu YR, Levy R, Ashby P, Tasker RR, Dostrovsky JO. Does stimulation of the GPi control dyskinesia by activating inhibitory axons?. *Mov Disord*. 2001;16(2):208–16. <https://doi.org/10.1002/mds.1046> PMID: 11295772
72. Hashimoto T, Elder CM, Okun MS, Patrick SK, Vitek JL. Stimulation of the subthalamic nucleus changes the firing pattern of pallidal neurons. *J Neurosci*. 2003;23(5):1916–23. <https://doi.org/10.1523/JNEUROSCI.23-05-01916.2003> PMID: 12629196
73. McCairn KW, Turner RS. Deep brain stimulation of the globus pallidus internus in the parkinsonian primate: local entrainment and suppression of low-frequency oscillations. *J Neurophysiol*. 2009;101(4):1941–60. <https://doi.org/10.1152/jn.91092.2008> PMID: 19164104
74. Doi T, Fan Y, Gold JL, Ding L. The caudate nucleus contributes causally to decisions that balance reward and uncertain visual information. *Elife*. 2020;9:e56694. <https://doi.org/10.7554/eLife.56694> PMID: 32568068
75. Wei W, Rubin JE, Wang X-J. Role of the indirect pathway of the basal ganglia in perceptual decision making. *J Neurosci*. 2015;35(9):4052–64. <https://doi.org/10.1523/JNEUROSCI.3611-14.2015> PMID: 25740532
76. Gurney KN, Humphries MD, Redgrave P. A new framework for cortico-striatal plasticity: behavioural theory meets in vitro data at the reinforcement-action interface. *PLoS Biol*. 2015;13(1):e1002034. <https://doi.org/10.1371/journal.pbio.1002034> PMID: 25562526
77. Baladron J, Nambu A, Hamker FH. The subthalamic nucleus-external globus pallidus loop biases exploratory decisions towards known alternatives: a neuro-computational study. *Eur J Neurosci*. 2019;49(6):754–67. <https://doi.org/10.1111/ejn.13666> PMID: 28833676
78. Fortin FA, Rainville FMD, Gardner MA, Parizeau M, Gagne C. DEAP: evolutionary algorithms made easy. *Journal of Machine Learning Research*. 2012;13(70):2171–5.