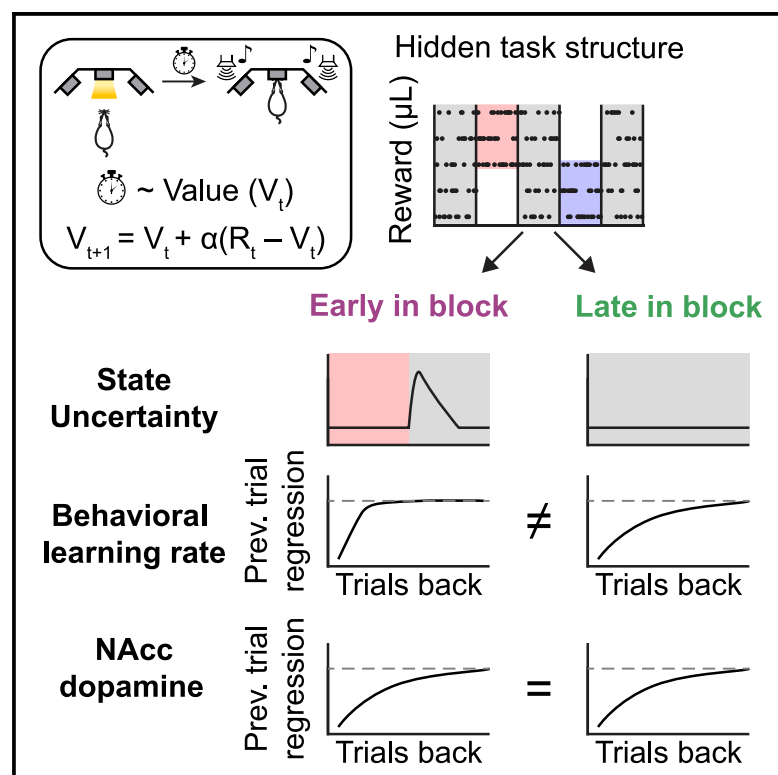


Dopamine transients encode reward prediction errors independent of learning rates

Graphical abstract



Authors

Andrew Mah, Carla E.M. Golden,
Christine M. Constantinople

Correspondence

constantinople@nyu.edu

In brief

Mah et al. find that rats use an uncertainty-based dynamic learning rate to adjust response vigor in a task with hidden states, which, they show, approximates normative Bayesian changepoint detection. However, despite strong behavioral effects, dopamine release in the nucleus accumbens does not reflect dynamic learning rates.

Highlights

- Rats use an uncertainty-based dynamic learning rate to modulate response vigor
- The dynamic learning rate approximates Bayesian changepoint detection models
- Dopamine release in the nucleus accumbens core does not differ by learning rate
- Accumbal dopamine release is modulated by temporal uncertainty about reward timing



Report

Dopamine transients encode reward prediction errors independent of learning rates

Andrew Mah,¹ Carla E.M. Golden,¹ and Christine M. Constantinople^{1,2,*}

¹Center for Neural Science, New York University, New York, NY, USA

²Lead contact

*Correspondence: constantinople@nyu.edu

<https://doi.org/10.1016/j.celrep.2024.114840>

SUMMARY

Biological accounts of reinforcement learning posit that dopamine encodes reward prediction errors (RPEs), which are multiplied by a learning rate to update state or action values. These values are thought to be represented by corticostriatal synaptic weights, which are updated by dopamine-dependent plasticity. This suggests that dopamine release reflects the product of the learning rate and RPE. Here, we characterize dopamine encoding of learning rates in the nucleus accumbens core (NAcc) in a volatile environment. Using a task with semi-observable states offering different rewards, we find that rats adjust how quickly they initiate trials across states using RPEs. Computational modeling and behavioral analyses show that learning rates are higher following state transitions and scale with trial-by-trial changes in beliefs about hidden states, approximating normative Bayesian strategies. Notably, dopamine release in the NAcc encodes RPEs independent of learning rates, suggesting that dopamine-independent mechanisms instantiate dynamic learning rates.

INTRODUCTION

Reinforcement learning describes how animals or agents learn the value of states and actions to select actions that maximize future expected rewards.¹ Reinforcement learning algorithms, including temporal-difference learning, update state and action values using reward prediction errors (RPEs), or the difference between experienced and expected rewards. The rate of error-driven learning is often assumed to be constant, but work across humans, monkeys, rats, and mice has found behavioral evidence for dynamic learning rates.^{2–8} In volatile environments, dynamic learning rates allow animals to learn faster when the world is changing and slower when the world is stable.^{9–11}

Dopaminergic inputs to the striatum, the input structure of the basal ganglia, are thought to convey a biological RPE. In reinforcement learning models of the basal ganglia, cortical inputs to the striatum convey the animal's state, with the strength of the synapse proportional to the expected future reward, or value, of that state.^{1,12,13} States with higher values have stronger synapses that are more likely to drive striatal action selection. The strengths of these corticostriatal synapses are updated via dopamine-dependent plasticity, proportional to dopamine RPEs.^{14–17} However, reinforcement learning algorithms update values proportional to the *product* of the RPE and learning rate. It is currently unclear whether dopamine conveys only the RPE, with other substrates dictating the learning rate, or whether dopamine encodes the RPE scaled by the learning rate. These scenarios are indistinguishable if learning rates are static. Here, we leveraged the fact that reinforcement learning is dynamic in chang-

ing environments to characterize dopamine encoding of learning rates in the nucleus accumbens core (NAcc) by recording dopamine release in rats performing a task with latent reward states.

RESULTS

Rats use a dynamic learning rate

We trained rats on a self-paced temporal wagering task with semi-observable reward blocks¹⁸ (Figures 1A and 1B). Rats were offered different volumes of water rewards (5, 10, 20, 40, and 80 μ L), cued by an auditory tone. On 75%–85% of trials, rewards were delivered after variable, unpredictable delays drawn from an exponential distribution. On 15%–25% of trials, rewards were withheld. The rats could choose to wait for the water reward or could opt out at any time to start a new trial. We introduced uncued blocks of trials with differing reward statistics; low blocks, which offered the three smallest rewards (5, 10, and 20 μ L), and high blocks, which offered the three largest rewards (20, 40, and 80 μ L), interleaved with mixed blocks, which offered all rewards (Figure 1B).

We measured the time between the rat's final poke in the reward or "opt-out" port and the start of the next trial (trial initiation time). Trial initiation times were inversely proportional to the value of the environment and provided a continuous behavioral readout of rats' estimates of state values.¹⁸ Rats were slower to initiate trials in low blocks compared to high blocks (Figures 1C and 1D). Furthermore, when we regressed initiation times against previous reward offers, coefficients were larger for more recent offers and gradually decreased for more distant trials (Figure 1E). This pattern is consistent with canonical



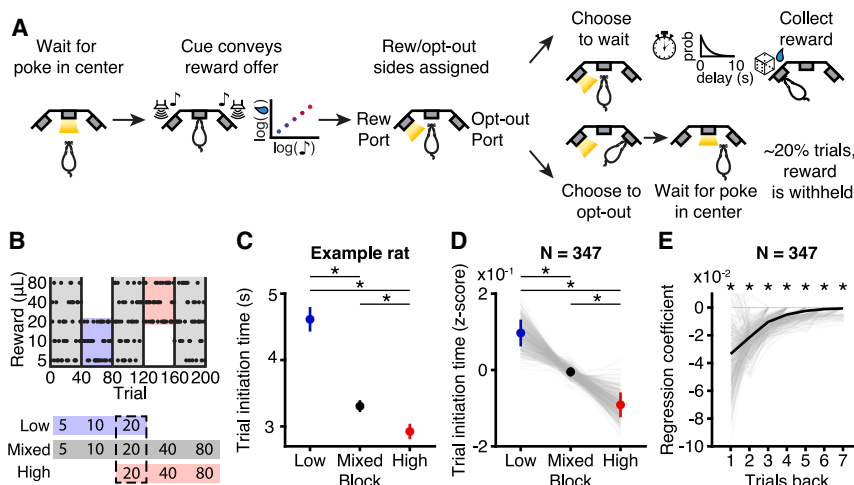


Figure 1. Trial initiation times are sensitive to previous rewards and blocks

(A) Task schematic. Rats initiate trials by poking in a center port. Auditory tones convey reward offers of different volumes, which are delivered with unpredictable delays, and are withheld on a subset of trials.

(B) Block structure.

(C and D) Average trial initiation times for a single rat ($p < 0.01$, Wilcoxon rank-sum test) (C) and all rats ($p \ll 0.001$, Wilcoxon signed-rank test, $N = 347$) (D) are sensitive to block.

(E) Regression coefficients of previous rewards predicting trial initiation times across rats ($p < 0.05$, Wilcoxon signed-rank test, $N = 347$). All error bars are mean \pm SEM.

reinforcement learning algorithms, which estimate the value of the environment as a recency-weighted average of previous rewards. We focused on trial initiation times because our previous work found that initiation times reflected state value estimates that were updated consistent with standard reinforcement learning algorithms, whereas willingness to wait during the reward delay reflected more complex state inference processes that are distinct from these standard algorithms.¹⁸

In reinforcement learning, the decay time of previous trial coefficients is directly proportional to the learning rate parameter, which is often assumed to be static. However, examination of trial initiation times aligned to transitions from low or high blocks into mixed blocks revealed two phases of learning: an initial phase of fast learning followed by slower dynamics later in the block, suggestive of higher learning rates immediately following block transitions (Figures 2A and S1). Previous work found that the overshoot in trial initiation times after block transitions could not be explained by a static learning rate¹⁸ and is robust across rats (Figure S2). Consistent with this result, when we regressed trial initiation times against previous rewards separately for the first and last 10 trials of each mixed block, rats integrated over fewer trials earlier in the block (i.e., had higher learning rates) compared to later in the block (Figures 2B, 2C, S3, and S4). Across rats, exponential functions fit to the regression coefficients had significantly smaller time constants (higher learning rates) for early versus late trials (Figure 2D).

We next separately fit a simple reinforcement learning model with a static learning rate to early and late trials (Figure 2E). The model estimated the value of the environment according to the recursive equation $V_{t+1} = V_t + \alpha(R_t - V_t)$, where $R_t - V_t$ is the RPE, and α is a learning rate that dictates the dynamics with which values are updated over trials. The trial initiation times were modeled as inversely proportional to this value, which captured the rat's behavior on held-out test data¹⁸ (Figure 2F). Models fit to early mixed-block trials had significantly higher learning rates (Figure 2G) and fewer significant previous trial coefficients (Figure 2H) compared to models fit to late mixed-block trials. Model comparison confirmed that learning rates fit to early

and late parts of the block provided better fits to held-out test data in early and late block trials, respectively (Figure 2H). Overall, these data suggest that rats use a dynamic learning rate that is higher following block transitions compared to later in the block.

Rats' dynamic learning rates reflect changing beliefs about reward blocks

Previous work found that animals adjust their learning rates depending on the volatility in the environment, since it is advantageous to learn faster in dynamic environments.^{3,5,7} To determine which features of the environment the rats were using to adjust their learning rates, we tested several classes of dynamic learning rate models (Figure 3A). Previous studies have proposed that subjects might scale their learning rates based on salient events or outcomes.^{19,20} We modeled salience as proportional to the log-reward offer on each trial, so larger rewards were assumed to be more salient than smaller rewards, which we call the Mackintosh model.¹⁹

Alternatively, previous work suggests that it is advantageous to modulate learning rates based on perceived volatility in the environment.^{3,5,21} When volatility is high, previous outcomes are less predictive of future states, so agents should use a higher learning rate to disregard distal outcomes. We tested two classes of volatility-based dynamic learning rate models. First, the model-free Pearce-Hall model adjusts the learning rate proportional to the unsigned RPE of the previous trial.²² Intuitively, large RPEs indicate that reward expectations are inaccurate, potentially because the environment is changing. Second, we characterized a "model-sensitive" model we developed previously,¹⁸ which we refer to as the Δ Belief model. This model uses Bayes' rule to infer the probability of being in each block given the rat's current and past reward offers and scales the learning rate by the trial-by-trial change in the belief about the block (i.e., change in posterior probability; STAR Methods).

We used these models to generate qualitative predictions about initiation time behavior. First, the models make distinct predictions about how the variance of the initiation times should change within mixed blocks. Higher learning rates imply higher

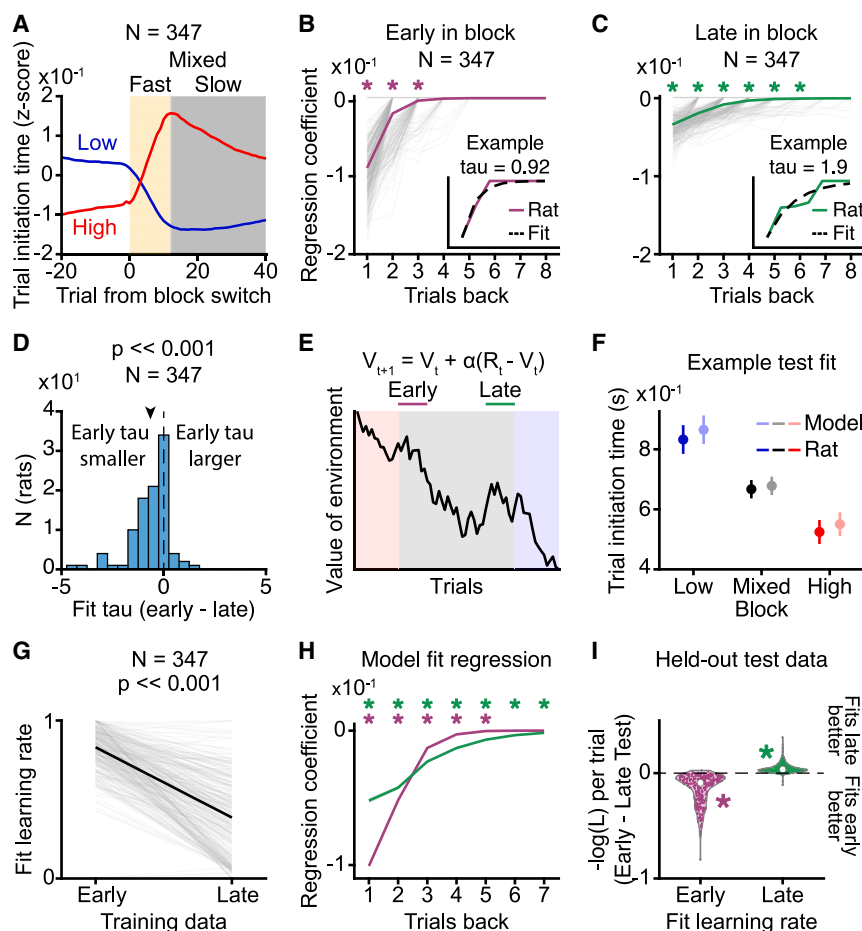


Figure 2. Rats use a higher learning rate early in mixed blocks

(A) Trial initiation times aligned to transitions into mixed blocks from low (blue) and high (red) blocks, smoothed with a causal filter of size 10 trials.

(B and C) Previous reward regression coefficients for (B) early and (C) late mixed-block trials (Wilcoxon signed-rank test, $N = 347$). Insets show regression coefficients from an example rat (solid) overlaid with exponential fits (dashed).

(D) Difference in time constant, τ , of exponential decay fit to early and late mixed-block regression coefficients across rats. An arrow indicates mean ($p < 0.001$, paired Wilcoxon signed-rank test, $N = 347$).

(E) Model schematic.

(F) Example model performance on held-out late-trial test data.

(G) Recovered learning rate parameters for early and late mixed-block training data across rats ($p < 0.001$, Wilcoxon signed-rank test, $N = 347$).

(H) Average previous trial regression from model fits to early (purple) or late (green) trials across rats ($N = 347$).

(I) Comparing negative log likelihood of early and late parameters on held-out test data. White circles indicate mean (Wilcoxon signed-rank test, $N = 347$). $*p < 0.05$. All error bars are mean \pm SEM.

variability in initiation times because the value estimate is updated more on each trial. Because the mixed blocks include all reward offers, the magnitude of RPEs is comparable both early and late in the blocks. Therefore, the Pearce-Hall model, which scales the learning rate by the previous unsigned RPE, predicts similar learning rates early and late in mixed blocks and, thus, similar initiation time variance (Figure 3C). Similarly, large rewards are equally likely early and late in mixed blocks, so the Mackintosh model, which scales the learning rate with the log-reward offer, predicts similar initiation time variance across mixed blocks (Figure 3C). However, the Δ Belief model, which scales the learning rate with the trial-by-trial change in the belief about the inferred block, predicts higher variance in initiation times early in the block, when beliefs are changing, compared to late, when beliefs are stable (Figures 3B and 3C). We measured the variance of the initiation times for the first and last 10 trials of each mixed block and computed their log ratio (negative values correspond to higher variance early in the block). Consistent with the Δ Belief model, but not the other models, the variance of rats' initiation times was higher in the early mixed block trials compared to later trials (Figure 3D).

Next, we explicitly tested a prediction from the Δ Belief model, that rats' trial-by-trial learning should depend on the change in belief about the block, controlling for RPEs. We focused on early

trials in the block, when there should be a broader range of beliefs to increase statistical power. We estimated RPEs using model parameters fit to behavior from the first 10 trials of each block. We compared the trial-by-trial change in initiation times for the same values of RPEs but conditioned on whether the Δ Belief on that trial was high or low ($>$ or $<$ 50th percentile). We found that rats changed their initiation times more for both positive and negative RPEs on trials with large changes in beliefs compared to trials with the same binned RPEs but small changes in beliefs (Figure 3E), consistent with the Δ Belief model (Figure 3F) but not the other dynamic learning rate models (Figures 3G and 3H). In summary, rats use knowledge about the underlying task structure to update their learning rates.

However, it remains unclear why rats would use the Δ Belief model. We found that the Δ Belief model approximates a normative algorithm known as Bayesian online changepoint detection.²³ Bayesian online changepoint detection uses sequential observations (e.g., reward offers) and a model of the environment to estimate the probability that the underlying distribution generating those observations (e.g., the reward block) has changed. These inferred transitions are called changepoints. Specifically, on each trial, the model generates a distribution over the number of trials since the last changepoint (run lengths) and chooses the run length that maximizes the posterior probability (Figures 3I–3K). A run length of K means the model estimates that a changepoint occurred K trials ago (Figures 3J and 3K). While Bayesian online changepoint detection can correctly estimate changepoints around block transitions in our task

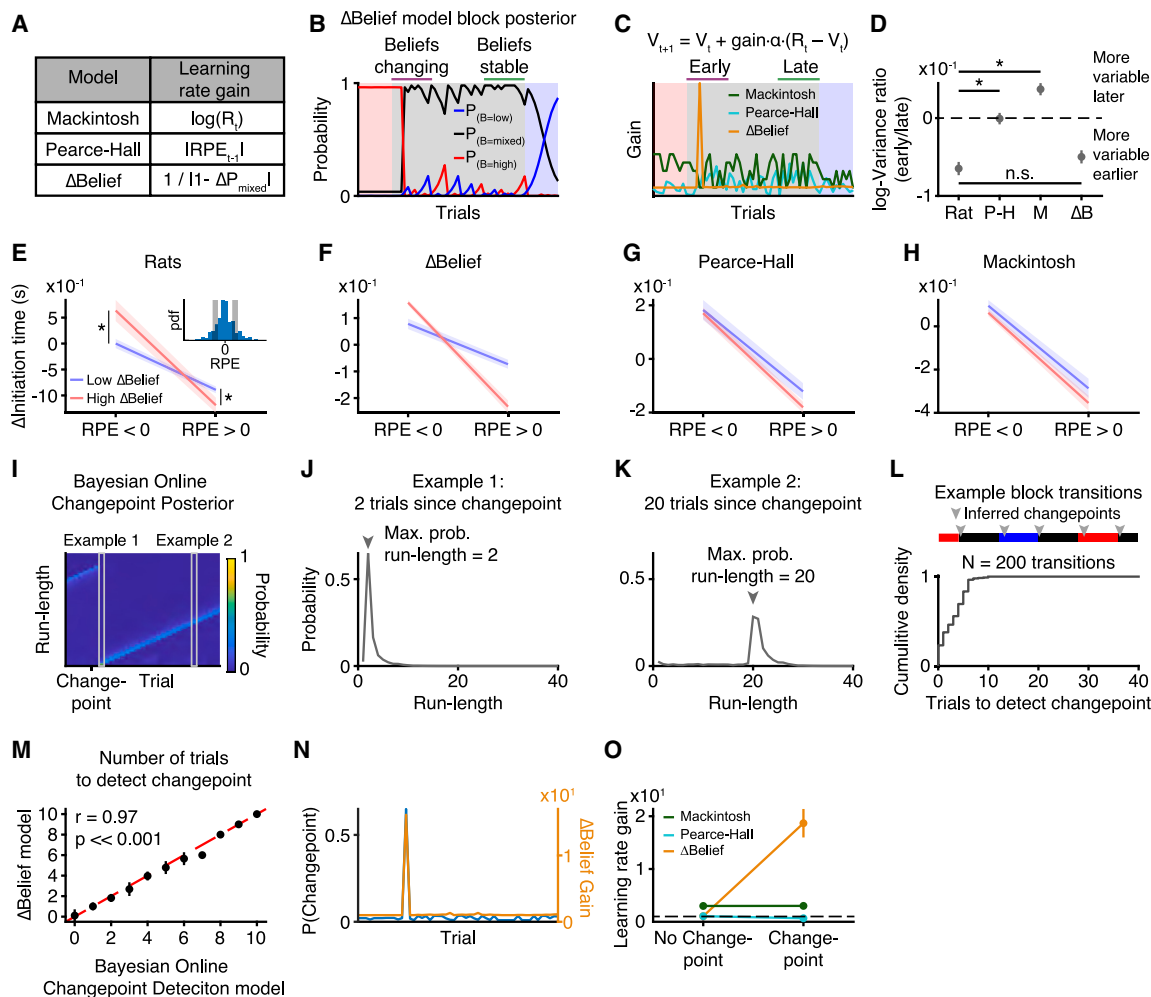


Figure 3. Rats use changing beliefs to modulate learning rates

(A) Explanation of models.
(B) Block posterior probability from the Δ Belief model. Beliefs change the most at block transitions.
(C) Examples of learning rate gain from (A).
(D) Trial initiation time variance early vs. late for dynamic learning rate models and rat data (Wilcoxon rank-sum test, $N = 347$).
(E) Change in trial initiation time for negative and positive RPE bins (inset) for large (red) or small (blue) changes in mixed-block belief ($>$ or $<$ 50th percentile, Wilcoxon signed-rank test, $N = 347$).
(F–H) Model predictions for the (F) Δ Belief, (G) Pearce-Hall, and (H) Mackintosh models.
(I) Bayesian online changepoint detection model run length posterior over trials around block transitions.
(J and K) Example changepoint posteriors for a trial (J) two trials from a changepoint or (K) 20 trials from a changepoint. Arrows indicate trials with maximum posterior probability.
(L) Simulated changepoint posterior for the Bayesian online changepoint detection model. Top: colored rectangles indicate true blocks, and triangles show inferred changepoints. Bottom: cumulative probability distribution for number of trials between block transitions and inferred changepoint ($N = 200$ simulated transitions).
(M) Average number of trials to detect the changepoint from the Δ Belief model and Bayesian online changepoint detection model ($N = 200$ simulated transitions; $p < 0.001$, Pearson correlation).
(N) Δ Belief model learning rate gain correlates with changepoint probability from the Bayesian online changepoint detection model.
(O) Learning rate gain for dynamic learning rate models on inferred changepoint trials. Asterisks indicate $p < 0.05$. All error bars are mean \pm SEM.

(Figure 3L), this model is computationally expensive. On each trial, the model iterates over all previous trials to find changepoints—the number of computations is a quadratic function of number of trials (STAR Methods). Strikingly, we found that the gain term from the Δ Belief model, which is less computationally expensive, is highly correlated with the changepoint probability

(Figure 3M). Over many simulated sessions, the Δ Belief gain was significantly higher on trials with inferred changepoints compared to other trials. There was no systematic relationship between the unsigned RPE (Pearce-Hall model) or log reward offer (Mackintosh model) and inferred changepoints (Figures 3N and 3O). The Δ Belief model can therefore provide a simple

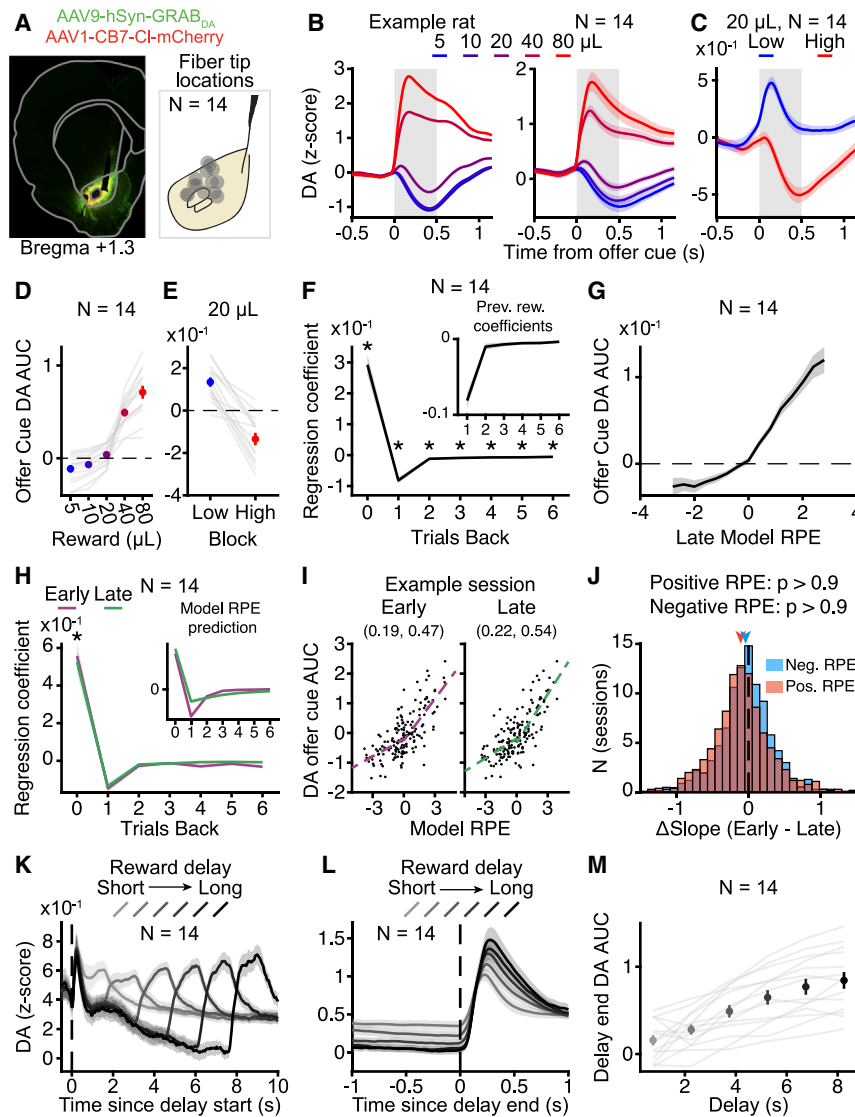


Figure 4. Dopamine release in the NAcc encodes RPEs independent of learning rates

(A) Example NAcc histology and recording site summary.

(B) NAcc dopamine aligned to offer cue by rewards in mixed blocks for an example rat (left) and averaged across rats (right, $N = 14$).

(C) Average NAcc dopamine response to offer cue for 20 μ L by block ($N = 14$).

(D) Average AUC from (B) (gray box, $N = 14$).

(E) Average AUC from (C) (gray box, $N = 14$).

(F) Reward history regression coefficients for NAcc dopamine signal. The inset shows previous reward coefficients (signed-rank test, $N = 14$).

(G) Offer cue AUC by binned late-model RPE ($N = 14$).

(H) Reward history regression coefficients for NAcc dopamine in early (purple) and late (green) trials. Paired signed-rank test, $N = 14$. The inset shows model prediction regression RPE vs. previous rewards for parameters fit to early vs. late trials.

(I) Example rat NAcc dopamine response by binned RPE for early (left) and late (right) trials (solid) with positive and negative RPE regressions (dashed). Numbers in parentheses indicate positive and negative slopes, respectively.

(J) Difference between early and late regression slopes for DA for positive (orange) and negative RPEs (blue) over individual sessions for all rats. Arrows indicate mean ($N = 994$ sessions, $p_{\text{early-plate}} > 0.9$, one-tailed Wilcoxon signed-rank test).

(K) NAcc dopamine aligned to the beginning of the delay period for rewarded mixed-block trials by reward delay ($N = 14$).

(L) NAcc dopamine aligned to the end of the delay period for rewarded mixed-block trials as a function of reward delay ($N = 14$).

(M) Baseline-corrected average AUC from (L) ($N = 14$). * $p < 0.05$. All error bars are mean \pm SEM.

approximation of a normative changepoint detection model that allows rats to detect changes in the environment and adjust their learning accordingly, similar to previous work that found behavioral evidence for simpler approximations of changepoint detection.²⁴

Dopamine release in the NAcc is not modulated by the dynamic learning rate

We next sought to find neural correlates of the dynamic learning rate in the NAcc, where dopamine is thought to mediate trial-by-trial learning by instantiating a biological RPE.^{25–38} Mechanistically, dopaminergic RPEs are thought to mediate plasticity at synapses onto medium spiny neurons, the principal cells of the striatum, to increase (or decrease) the likelihood of taking actions in a state that previously produced positive (or negative) RPEs. This account of dopamine function implicitly assumes that dopamine release represents the product of the learning rate and the RPE and so should reflect dynamic learning rates in our task.

We focused our recordings on the NAcc, which is thought to determine the vigor of motivated behaviors.³⁹ Recent work from our lab has also found that dopamine RPEs in the NAcc causally determine initiation times on subsequent trials.⁴⁰ We recorded dopamine release in the NAcc ($N = 14$; Figure S4) using fiber photometry and a fluorescent G-protein-coupled receptor-activation-based DA sensor (GRAB_{DA}).⁴¹ We observed robust phasic dopamine responses at the time of the offer cue that were consistent with an RPE. First, an RPE should correlate with reward offer. We found that NAcc dopamine release scaled monotonically with offered reward volume in mixed blocks, with dips for smaller rewards (Figures 4B and 4D). An RPE signal should also scale inversely with expectations. Focusing on 20- μ L trials, which appeared in all blocks, dopamine increased on 20- μ L trials in low blocks, when reward expectations were low, and decreased on 20- μ L trials in high blocks, when expectations were high (Figures 4C and 4E). Third, when we regressed NAcc dopamine release (area under the curve from 0 to 0.5 s; AUC) against reward history, we found positive coefficients for the current trial and negative coefficients for previous trials, a

hallmark of RPE encoding²⁹ (Figure 4F). Finally, NAcc dopamine release correlated with RPEs estimated from the behavioral model (Figures 4G and 4I). We related NAcc dopamine from the first and last 10 trials of mixed blocks to RPEs estimated from separate model fits to early and late block trials, respectively.

Next, we compared NAcc dopamine release early vs. late in the block when the rats used different learning rates. If dopamine release is the product of the RPE and the learning rate, then dopamine should be influenced by fewer trials in the past when the learning rate is higher (Figure 4H inset). However, there was no significant difference between previous trial regression coefficients fit to NAcc dopamine release early or late in the block (Figure 4H). This hypothesis also predicts that dopamine release would be greater for the same RPE earlier in the block, when the learning rate is higher, compared to later in the block, when beliefs are stable and learning rates are lower. To account for nonlinear encoding of RPEs, we fit separate regression lines to dopamine encoding (AUC) of positive and negative model RPEs for each session. There was no significant difference in the slope parameters fit to early or late trials in each session across rats, and this was true for both positive and negative RPEs (Figures 4I, 4J, and S5). Therefore, despite strong behavioral evidence for higher learning rates earlier in mixed blocks when beliefs about hidden states are changing, we did not find differential NAcc dopamine dynamics between early and late block trials.

Previous work found that the activity of ventral tegmental area dopamine neurons can reflect beliefs about hidden task states if rewards are probabilistic and variable in their timing.⁴² We therefore examined NAcc dopamine during the delay period. Because rewards were omitted on a subset of trials, as animals waited for the reward, there is ambiguity about whether they are in a rewarded or unrewarded trial. Sorting rewarded trials based on delay duration revealed negative dopamine ramps during the delay (Figure 4K), consistent with previous studies.⁴² These negative ramps were also apparent as different baselines when the data were aligned to the cue indicating that a reward was available (Figure 4L). These ramps have been interpreted as moment-by-moment negative RPEs as the animals wait without receiving a reward.³² When a cue indicated that a reward was available, there was a large phasic dopamine response in the NAcc whose magnitude scaled with the delay (Figures 4K–4M). This pattern cannot be captured by traditional model-free temporal difference learning and requires additional knowledge about the task structure⁴²; the probability of being in an unrewarded trial increases over time, so the reward cue is more unexpected given beliefs about the trial. Therefore, NAcc dopamine reflects beliefs about hidden states at the time of probabilistic rewards but not about latent reward blocks at the time of the offer cue.

DISCUSSION

Our results offer a nuanced, “model-sensitive” view of dopamine activity in which different aspects of state inference differentially modulate dopamine release in the NAcc. On one hand, NAcc dopamine at offer cues encodes RPEs independent of

inference about the latent reward block despite beliefs about the blocks strongly influencing reinforcement learning and behavior (Figures 4H–4J). This is consistent with traditional model-free interpretations of dopamine activity. However, at the end of the delay period, we found delay-dependent patterns of NAcc dopamine release that cannot be explained by model-free temporal difference learning, consistent with previous work⁴² (Figures 4K–4M). These results add to a body of work expanding dopamine beyond traditional, model-free learning algorithms.^{35,43}

It remains unclear why some aspects of state inference modulate NAcc dopamine while others do not. In a conditioning paradigm with probabilistic rewards, during the delay period, there is ambiguity about the state the animal is in: whether it is in the delay period or whether it has transitioned to the inter-trial interval on an unrewarded trial. In other words, there is *uncertainty* about the current state. However, this type of uncertainty is qualitatively different from the uncertainty about the latent reward blocks in our task. Uncertainty in the timing of probabilistic rewards reflects *expected*, or irreducible, uncertainty.²¹ Even with a perfect model of the world, there is inherent uncertainty in the outcome of stochastic events—one can know the probabilities of a coin toss but cannot predict the outcome of the next flip. This contrasts with *unexpected* uncertainty, which is related to the volatility of the environment—a perfect model of the world now may not be perfect if the world changes.^{21,44} In other words, unexpected uncertainty about the reward blocks can be reduced with additional observations, while expected uncertainty about probabilistic reward timing is irreducible. We find that dopamine release in the NAcc reflects expected but not unexpected uncertainty. Together, these results suggest that different types of uncertainty may map onto neurobiologically distinct mechanisms with potentially dissociable consequences for learning.^{8,21,34,44–46}

Previous work has suggested that dopamine release in the NAcc directly encodes the learning rate,⁴⁷ implicating dopamine in a broader class of policy-learning algorithms beyond traditional value learning. However, our work examines distinct phases of learning from Coddington et al.⁴⁷ We exclusively recorded from expert, as opposed to naive, animals. Presumably, our expert rats have already learned their final behavioral policy and must only learn about the current state of the environment. By contrast, task-naïve animals are simultaneously learning the associative structure of the environment and how to optimally behave in that environment. Such distinct learning goals likely engage NAcc dopamine differently and could explain the differences in our findings.

One key future question is what could be driving the dynamic learning rate at the level of synaptic plasticity. Previous work has shown that plasticity at corticostriatal synapses depends on the coordinated activity of dopamine^{14–17} with other neuromodulators, like acetylcholine⁴⁸ and serotonin.^{49–51} Serotonin neurons, which project to the NAcc from the dorsal raphe nucleus,⁵² have been shown to encode unexpected uncertainty⁷ and can causally influence learning rates in mice.⁵³ Other neuromodulators that have been hypothesized to encode unexpected uncertainty, like norepinephrine,⁴⁴ do not strongly project to the NAcc^{54–56} but could influence learning in other neural circuits.

Future studies clarifying task-related dynamics of other neuro-modulators may elucidate the circuit mechanisms that combine dopaminergic RPEs with trial-by-trial changes in beliefs to modulate the rate of learning at behavioral and synaptic levels.

Limitations of the study

Our current study focused on recording dopamine release in the NAcc. However, given recent findings that show considerable heterogeneity in dopamine activity across the striatum,^{57–68} additional studies are required to understand how dopamine activity across striatal subregions is modulated by hidden-state inference. More specific recording techniques, such as optogenetically tagged recordings or cell-type-specific fluorescent imaging, may help elucidate the dynamics of different classes of dopamine neurons and their implications for behavior.

RESOURCE AVAILABILITY

Lead contact

Requests for further information, resources, and reagents should be directed to and will be fulfilled by the lead contact, Christine Constantinople (cmc472@nyu.edu).

Materials availability

This study did not generate new unique reagents.

Data and code availability

- Behavioral and photometry data have been deposited at Zenodo and are publicly available as of the date of publication. Accession numbers are listed in the [key resources table](#).
- All original code has been deposited at GitHub and is publicly available as of the date of publication. DOIs are listed in the [key resources table](#).
- Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

ACKNOWLEDGMENTS

This work was supported by the National Institutes of Health via grant R00MH111926 (to C.M.C.), grant DP2MH126376 (to C.M.C.), grant F32MH125448 (to C.E.M.G.), grant 5T32MH019524 (to C.E.M.G. and A.M.), grant 5T90DA043219 (to A.M.), and grant F31MH130121 (to A.M.); an Alfred P. Sloan Research Fellowship (to C.M.C.); a Klingenstein-Simons Fellowship in Neuroscience (to C.M.C.); a McKnight Scholars Award (to C.M.C.); an NSF CAREER Award (to C.M.C.); and the Simons Foundation via grant 855332 (to C.E.M.G.). We thank members of the Constantinople lab for feedback and helpful discussions. We thank research technicians in the Constantinople lab for animal training.

AUTHOR CONTRIBUTIONS

C.M.C. and A.M. designed the study. C.E.G. collected the photometry data and contributed to photometry analysis. A.M. developed the behavioral model and performed behavioral and photometry analyses. A.M. prepared the figures. C.M.C. and A.M. wrote the manuscript. C.M.C. supervised the project.

DECLARATION OF INTERESTS

The authors declare no competing interests.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS**
- **METHOD DETAILS**
 - Behavioral training
 - Training for male and female rats
 - Criteria for including behavioral data
 - Behavioral modeling
 - Dynamic learning rate models
 - Fitting and evaluating models
 - Bayesian Online Change-point Detection model
 - Sterotaxic surgeries
 - Photometry
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
 - Sensitivity to reward blocks
 - Block transition dynamics
 - Previous reward regression
 - Photometry

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.celrep.2024.114840>.

Received: April 15, 2024

Revised: August 19, 2024

Accepted: September 20, 2024

REFERENCES

1. Sutton, R.S., and Barto, A.G. (2018). *Reinforcement Learning: An Introduction* (MIT press).
2. Nassar, M.R., Wilson, R.C., Heasly, B., and Gold, J.I. (2010). An approximately bayesian delta-rule model explains the dynamics of belief updating in a changing environment. *J. Neurosci.* 30, 12366–12378.
3. Behrens, T.E.J., Woolrich, M.W., Walton, M.E., and Rushworth, M.F.S. (2007). Learning the value of information in an uncertain world. *Nat. Neurosci.* 10, 1214–1221.
4. Hayden, B.Y., Heilbronner, S.R., Pearson, J.M., and Platt, M.L. (2011). Surprise signals in anterior cingulate cortex: neuronal encoding of unsigned reward prediction errors driving adjustment in behavior. *J. Neurosci.* 31, 4178–4187.
5. Nassar, M.R., Rumsey, K.M., Wilson, R.C., Parikh, K., Heasly, B., and Gold, J.I. (2012). Rational regulation of learning dynamics by pupil-linked arousal systems. *Nat. Neurosci.* 15, 1040–1046.
6. Funamizu, A., Ito, M., Doya, K., Kanzaki, R., and Takahashi, H. (2012). Uncertainty in action-value estimation affects both action choice and learning rate of the choice behaviors of rats. *Eur. J. Neurosci.* 35, 1180–1189.
7. Grossman, C.D., Bari, B.A., and Cohen, J.Y. (2022). Serotonin neurons modulate learning rate through uncertainty. *Curr. Biol.* 32, 586–599.e7.
8. McGuire, J.T., Nassar, M.R., Gold, J.I., and Kable, J.W. (2014). Functionally dissociable influences on learning rate in a dynamic environment. *Neuron* 84, 870–881.
9. Amari, S. (1967). A theory of adaptive pattern classifiers. *IEEE Trans. Electron. Comput.* EC-16, 299–307.
10. Sutton, R.S. (1992). Adapting bias by gradient descent: An incremental version of delta-bar-delta. In *AAAI*, 92 (Citeseer), pp. 171–176.
11. Murata, N., Kawanabe, M., Ziehe, A., Müller, K.-R., and Amari, S.-i. (2002). On-line learning in changing environments with applications in supervised and unsupervised learning. *Neural Network.* 15, 743–760.
12. Joel, D., Niv, Y., and Ruppel, E. (2002). Actor-critic models of the basal ganglia: New anatomical and computational perspectives. *Neural Network.* 15, 535–547.

13. Doya, K. (2007). Reinforcement learning: Computational theory and biological mechanisms. *HFSP J.* 1, 30–40.
14. Centonze, D., Gubellini, P., Picconi, B., Calabresi, P., Giacomini, P., and Bernardi, G. (1999). Unilateral dopamine denervation blocks corticostriatal ltp. *J. Neurophysiol.* 82, 3575–3579.
15. Kerr, J.N., and Wickens, J.R. (2001). Dopamine d-1/d-5 receptor activation is required for long-term potentiation in the rat neostriatum in vitro. *J. Neurophysiol.* 85, 117–124.
16. Reynolds, J.N., Hyland, B.I., and Wickens, J.R. (2001). A cellular mechanism of reward-related learning. *Nature* 413, 67–70.
17. Shen, W., Flajolet, M., Greengard, P., and Surmeier, D.J. (2008). Dichotomous dopaminergic control of striatal synaptic plasticity. *Science* 321, 848–851.
18. Mah, A., Schierack, S.S., Bossio, V., and Constantinople, C.M. (2023). Distinct value computations support rapid sequential decisions. *Nat. Commun.* 14, 7573.
19. Mackintosh, N.J. (1975). A theory of attention: Variations in the associability of stimuli with reinforcement. *Psychol. Rev.* 82, 276–298.
20. Iigaya, K. (2016). Adaptive learning and decision-making under uncertainty by metaplastic synapses guided by a surprise detection system. *Elife* 5, e18073.
21. Soltani, A., and Izquierdo, A. (2019). Adaptive learning under expected and unexpected uncertainty. *Nat. Rev. Neurosci.* 20, 635–644.
22. Pearce, J.M., and Hall, G. (1980). A model for pavlovian learning: variations in the effectiveness of conditioned but not of unconditioned stimuli. *Psychol. Rev.* 87, 532–552.
23. Adams, R.P., and MacKay, D.J. (2007). Bayesian online changepoint detection. Preprint at arXiv. <https://doi.org/10.48550/arXiv:0710.3742>.
24. Wilson, R.C., Nassar, M.R., and Gold, J.I. (2013). A mixture of delta-rules approximation to bayesian inference in change-point problems. *PLoS Comput. Biol.* 9, e1003150.
25. Olds, J. (1958). Self-stimulation of the brain: Its use to study local effects of hunger, sex, and drugs. *Science* 127, 315–324.
26. Corbett, D., and Wise, R.A. (1980). Intracranial self-stimulation in relation to the ascending dopaminergic systems of the midbrain: a moveable electrode mapping study. *Brain Res.* 185, 1–15.
27. Schultz, W., Dayan, P., and Montague, P.R. (1997). A neural substrate of prediction and reward. *Science* 275, 1593–1599.
28. Waelti, P., Dickinson, A., and Schultz, W. (2001). Dopamine responses comply with basic assumptions of formal learning theory. *Nature* 412, 43–48.
29. Bayer, H.M., and Glimcher, P.W. (2005). Midbrain dopamine neurons encode a quantitative reward prediction error signal. *Neuron* 47, 129–141.
30. Cohen, J.Y., Haesler, S., Vong, L., Lowell, B.B., and Uchida, N. (2012). Neuron-type-specific signals for reward and punishment in the ventral tegmental area. *Nature* 482, 85–88.
31. Day, J.J., Roitman, M.F., Wightman, R.M., and Carelli, R.M. (2007). Associative learning mediates dynamic shifts in dopamine signaling in the nucleus accumbens. *Nat. Neurosci.* 10, 1020–1028.
32. Kim, H.R., Malik, A.N., Mikhael, J.G., Bech, P., Tsutsui-Kimura, I., Sun, F., Zhang, Y., Li, Y., Watabe-Uchida, M., Gershman, S.J., and Uchida, N. (2020). A unified framework for dopamine signals across timescales. *Cell* 183, 1600–1616.e25.
33. Steinberg, E.E., Keiflin, R., Boivin, J.R., Witten, I.B., Deisseroth, K., and Janak, P.H. (2013). A causal link between prediction errors, dopamine neurons and learning. *Nat. Neurosci.* 16, 966–973.
34. Parker, N.F., Cameron, C.M., Taliaferro, J.P., Lee, J., Choi, J.Y., Davidson, T.J., Daw, N.D., and Witten, I.B. (2016). Reward and choice encoding in terminals of midbrain dopamine neurons depends on striatal target. *Nat. Neurosci.* 19, 845–854.
35. Sharpe, M.J., Chang, C.Y., Liu, M.A., Batchelor, H.M., Mueller, L.E., Jones, J.L., Niv, Y., and Schoenbaum, G. (2017). Dopamine transients are sufficient and necessary for acquisition of model-based associations. *Nat. Neurosci.* 20, 735–742.
36. Tsai, H.-C., Zhang, F., Adamantidis, A., Stuber, G.D., Bonci, A., De Lecea, L., and Deisseroth, K. (2009). Phasic firing in dopaminergic neurons is sufficient for behavioral conditioning. *Science* 324, 1080–1084.
37. Hamid, A.A., Pettibone, J.R., Mabrouk, O.S., Hetrick, V.L., Schmidt, R., Vander Weele, C.M., Kennedy, R.T., Aragona, B.J., and Berke, J.D. (2016). Mesolimbic dopamine signals the value of work. *Nat. Neurosci.* 19, 117–126.
38. Adamantidis, A.R., Tsai, H.-C., Boutrel, B., Zhang, F., Stuber, G.D., Budygin, E.A., Touriño, C., Bonci, A., Deisseroth, K., and de Lecea, L. (2011). Optogenetic interrogation of dopaminergic modulation of the multiple phases of reward-seeking behavior. *J. Neurosci.* 31, 10829–10835.
39. Floresco, S.B. (2015). The nucleus accumbens: an interface between cognition, emotion, and action. *Annu. Rev. Psychol.* 66, 25–52.
40. Golden, C.E.M., Kaur, D., Mah, A., Martin, A.C., Levy, D.H., Yamaguchi, T., Lin, D., Aoki, C., and Constantinople, C.M. (2023). Estrogenic control of reward prediction errors and reinforcement learning. Preprint at bioRxiv. <https://doi.org/10.1101/2023.12.09.570945-12>.
41. Sun, F., Zhou, J., Dai, B., Qian, T., Zeng, J., Li, X., Zhuo, Y., Zhang, Y., Wang, Y., Qian, C., et al. (2020). Next-generation grab sensors for monitoring dopaminergic activity in vivo. *Nat. Methods* 17, 1156–1166.
42. Starkweather, C.K., Babayan, B.M., Uchida, N., and Gershman, S.J. (2017). Dopamine reward prediction errors reflect hidden-state inference across time. *Nat. Neurosci.* 20, 581–589.
43. Gershman, S.J., and Uchida, N. (2019). Believing in dopamine. *Nat. Rev. Neurosci.* 20, 703–714.
44. Angela, J.Y., and Dayan, P. (2005). Uncertainty, neuromodulation, and attention. *Neuron* 46, 681–692.
45. O'reilly, J.X. (2013). Making predictions in a changing world—inference, uncertainty, and learning. *Front. Neurosci.* 7, 33773.
46. Payzan-LeNestour, E., and Bossaerts, P. (2011). Risk, unexpected uncertainty, and estimation uncertainty: Bayesian learning in unstable settings. *PLoS Comput. Biol.* 7, e1001048.
47. Coddington, L.T., Lindo, S.E., and Dudman, J.T. (2023). Mesolimbic dopamine adapts the rate of learning from action. *Nature* 614, 294–302.
48. Reynolds, J.N.J., Avvisati, R., Dodson, P.D., Fisher, S.D., Oswald, M.J., Wickens, J.R., and Zhang, Y.-F. (2022). Coincidence of cholinergic pauses, dopaminergic activation and depolarisation of spiny projection neurons drives synaptic plasticity in the striatum. *Nat. Commun.* 13, 1296.
49. Burke, D.A., and Alvarez, V.A. (2022). Serotonin receptors contribute to dopamine depression of lateral inhibition in the nucleus accumbens. *Cell Rep.* 39, 110795.
50. Pommer, S., Akamine, Y., Schiffmann, S.N., de Kerchove d'Exaerde, A., and Wickens, J.R. (2021). The effect of serotonin receptor 5-HT1b on lateral inhibition between spiny projection neurons in the mouse striatum. *J. Neurosci.* 41, 7831–7847.
51. Mathur, B.N., Capik, N.A., Alvarez, V.A., and Lovinger, D.M. (2011). Serotonin induces long-term depression at corticostriatal synapses. *J. Neurosci.* 31, 7402–7411.
52. Van Bockstaele, E.J., Biswas, A., and Pickel, V.M. (1993). Topography of serotonin neurons in the dorsal raphe nucleus that send axon collaterals to the rat prefrontal cortex and nucleus accumbens. *Brain Res.* 624, 188–198.
53. Iigaya, K., Fonseca, M.S., Murakami, M., Mainen, Z.F., and Dayan, P. (2018). An effect of serotonergic stimulation on learning rates for rewards apparent after long intertrial intervals. *Nat. Commun.* 9, 2477.
54. Allin, R., Russell, V.A., Lamm, M.C., and Taljaard, J.J. (1988). Regional distribution of monoamines in the nucleus accumbens of the rat. *Neurochem. Res.* 13, 937–942.
55. Delfs, J.M., Zhu, Y., Druhan, J.P., and Aston-Jones, G.S. (1998). Origin of noradrenergic afferents to the shell subregion of the nucleus

- accumbens: anterograde and retrograde tract-tracing studies in the rat. *Brain Res.* 806, 127–140.
56. McKittrick, C.R., and Abercrombie, E.D. (2007). Catecholamine mapping within nucleus accumbens: differences in basal and amphetamine-stimulated efflux of norepinephrine and dopamine in shell and core. *J. Neurochem.* 100, 1247–1256.
 57. Elum, J.E., Szelenyi, E.R., Juarez, B., Murry, A.D., Loginov, G., Zamorano, C.A., Gao, P., Wu, G., Ng-Evans, S., Xu, X., et al. (2024). Distinct dynamics and intrinsic properties in ventral tegmental area populations mediate reward association and motivation. Preprint at bioRxiv. <https://doi.org/10.1101/2024.02.05.578997>.
 58. Engelhard, B., Finkelstein, J., Cox, J., Fleming, W., Jang, H.J., Ornelas, S., Koay, S.A., Thiberge, S.Y., Daw, N.D., Tank, D.W., and Witten, I.B. (2019). Specialized coding of sensory, motor and cognitive variables in vta dopamine neurons. *Nature* 570, 509–513.
 59. Howe, M.W., and Dombeck, D.A. (2016). Rapid signalling in distinct dopaminergic axons during locomotion and reward. *Nature* 535, 505–510.
 60. Heymann, G., Jo, Y.S., Reichard, K.L., McFarland, N., Chavkin, C., Palmiter, R.D., Soden, M.E., and Zweifel, L.S. (2020). Synergy of distinct dopamine projection populations in behavioral reinforcement. *Neuron* 105, 909–920.e5.
 61. Collins, A.L., and Saunders, B.T. (2020). Heterogeneity in striatal dopamine circuits: Form and function in dynamic reward seeking. *J. Neurosci. Res.* 98, 1046–1069.
 62. Lammel, S., Ion, D.I., Roeper, J., and Malenka, R.C. (2011). Projection-specific modulation of dopamine neuron synapses by aversive and rewarding stimuli. *Neuron* 70, 855–862.
 63. Horvitz, J.C. (2000). Mesolimbocortical and nigrostriatal dopamine responses to salient non-reward events. *Neuroscience* 96, 651–656.
 64. Cai, L.X., Pizano, K., Gundersen, G.W., Hayes, C.L., Fleming, W.T., Holt, S., Cox, J.M., and Witten, I.B. (2020). Distinct signals in medial and lateral vta dopamine neurons modulate fear extinction at different times. *Elife* 9, e54936.
 65. de Jong, J.W., Liang, Y., Verharen, J.P.H., Fraser, K.M., and Lammel, S. (2024). State and rate-of-change encoding in parallel mesoaccumbal dopamine pathways. *Nat. Neurosci.* 27, 309–318.
 66. Saunders, B.T., Richard, J.M., Margolis, E.B., and Janak, P.H. (2018). Dopamine neurons create pavlovian conditioned stimuli with circuit-defined motivational properties. *Nat. Neurosci.* 21, 1072–1083.
 67. Brischoux, F., Chakraborty, S., Brierley, D.I., and Ungless, M.A. (2009). Phasic excitation of dopamine neurons in ventral vta by noxious stimuli. *Proc. Natl. Acad. Sci. USA* 106, 4894–4899.
 68. Badrinarayan, A., Wescott, S.A., Vander Weele, C.M., Saunders, B.T., Couturier, B.E., Maren, S., and Aragona, B.J. (2012). Aversive stimuli differentially modulate real-time dopamine transmission dynamics within the nucleus accumbens core and shell. *J. Neurosci.* 32, 15779–15790.
 69. Niv, Y., Daw, N.D., Joel, D., and Dayan, P. (2007). Tonic dopamine: opportunity costs and the control of response vigor. *Psychopharmacology* 191, 507–520.
 70. Creamer, M.S., Chen, K.S., Leifer, A.M., and Pillow, J.W. (2022). Correcting motion induced fluorescence artifacts in two-channel neural imaging. *PLoS Comput. Biol.* 18, e1010421.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Antibodies		
GFP Polyclonal Antibody	Thermo Fisher Scientific	Cat# A11122; RRID: AB_221569
Goat anti-Rabbit IgG (H + L) Cross-Adsorbed Secondary Antibody, Alexa Fluor™ 488	Thermo Fisher Scientific	Cat# A11008; RRID: AB_143165
Bacterial and virus strains		
pAAV-hsyn-GRAB_DA2h	Sun et al. ⁴¹	Addgene AAV9; 140554-AAV9
pENN.AAV.CB7.Cl.mCherry.WPRE.RBG (AAV9)	James M. Wilson	Addgene AAV9; 105544-AAV9
Deposited data		
Raw and analyzed behavioral data	This paper	https://doi.org/10.5281/zenodo.13748709
Raw and analyzed photometry data	Golden et al. ⁴⁰	https://doi.org/10.5281/zenodo.13891951
Experimental models: Organisms/strains		
Long-Evans Rats	Hilltop Lab Animals	Hla®(LE)CVF®
Long-Evans Rat	Charles River	006
Software and algorithms		
MATLAB	MathWorks	R2023a, R2024a
Custom analysis for behavioral and photometry data	This paper	https://doi.org/10.5281/zenodo.13819804

EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS

A total of 347 Long-evans rats (*Rattus norvegicus*; 215 males, 132 females) between the ages of 6 and 24 months old. We previously found no differences between male and female rats.¹⁸ We found no effect of age on the main behavioral findings (Figure S6). This cohort included 24 TH-Cre rats, 8 ADORA2A-Cre, and 3 DRD1-Cre rats. We also did not find any effect of genotype on the main behavioral findings (Figure S7). Animal procedures were approved by the New York University Animal Welfare Committee (UAWC #2021-1120) and carried out in accordance with National Institute of Health standards.

Rats were typically pair-housed. To motivate behavioral performance, rats were water restricted from Monday to Friday, during which time they received water during behavioral training sessions, typically 90 min, followed by 20 min of *ad libitum* water. Rats were given *ad libitum* water following training on Friday through mid-day Sunday. Rats were weighed daily.

METHOD DETAILS

Behavioral training

We have previously published a detailed description of the behavioral shaping procedure for this task.¹⁸ Rats performed a self-paced temporal wagering task. Rats initiated trials by maintaining a nose poke in the center port for a variable period drawn from a uniform distribution over [0.8, 1.2] seconds. As the rat maintained the nose poke, the reward offer on that trial was conveyed by an auditory tone [1, 2, 4, 8, 16 kHz], which mapped onto one of five rewards ([5, 10, 20, 40, 80μL] for males, [4, 8, 16, 32, 64μL] for females). Following the reward tone presentation, rats could either wait a random delay drawn from an exponential distribution with mean of 2.5 s to receive their reward, or could opt-out at any time to immediately start a new trial. On 15–25% of trials (catch trials), reward was withheld to force the rats to exercise the opt-out option.

Training for male and female rats

We collected data from both male and female rats (215 males, 134 females). Male and female rats were trained with the same shaping procedure. Early cohorts of female rats experienced the same reward set as the males. However, because female rats are smaller, they consumed less water and performed substantially fewer trials than the males. Therefore, to obtain sufficient behavioral trials from both, females reward offers were slightly reduced while maintaining the logarithmic spacing: [4, 8, 16, 32, 64μL]. For behavioral analysis, reward volumes were treated as equivalent to the corresponding volume for the male rats (e.g., 16μL trials for female rats were treated the same as 20μL trials for male rats). The auditory tones were identical to those used for male rats. We did not observe any significant differences between the male and female rats, in terms of contextual effects, or behavioral dynamics at block transitions.¹⁸ Photometry data in this study was collected from females.

Criteria for including behavioral data

To determine whether rats sufficiently understood the mapping between auditory cues and water reward volumes, we evaluated their wait times on catch trials as a function of offered rewards. For each session, we first removed wait times that were greater than two standard deviations from the mean, which likely reflected lapses in attention/task disengagement. Next, we regressed wait time against offered reward. We included sessions with significant positive slopes that preceded at least one other session with a positive slope. We excluded trials with trial initiation times above the 99th percentile of the rat's cumulative trial initiation time distribution pooled over sessions.

Behavioral modeling

To model trial initiation times, we developed computational models based on,^{18,69} which describe the optimal trial initiation time, TI , given the value of the environment, V , as

$$TI = \frac{D}{V}$$

where D is a scale parameter. We developed multiple computational models that instantiated different algorithms for estimating the value of the environment. We estimate the value of the environment on trial t , V_t , by the recursive formula

$$V_{t+1} = V_t + \alpha_t(R_t - V_t)$$

where $\alpha_t = g_t \cdot \alpha_0$ is the learning rate, g is the learning rate gain, and R_t is the \log_2 (reward) on the current trial. For the static learning rate model, $g = 1$ for all trials.

Dynamic learning rate models

We tested several models of dynamic learning rates.

- (1) **Mackintosh surprise model** In this model, the gain on the learning rate is proportional to the salience of that trial,¹⁹ which we assumed to be directly proportional to the reward offer volume on that trials so,

$$g_t = \log_2(R_t)$$

where α_0 is the base learning rate.

- (2) **Pearce-Hall model** In this model, the learning rate gain is directly proportional to the inferred volatility of the environment. Volatility in this model is “model-free” as estimated as the unsigned RPE on the previous trial,²² so

$$g_t = |\text{RPE}_{t-1}|$$

where RPE_{t-1} is the reward prediction error on the previous trial.

- (3) **Δ Belief model** In this model, as in (2), the learning rate is directly proportional to the inferred volatility of the environment. In this model, volatility is calculated using the trial-by-trial change in the belief of being in a mixed block, using Bayes rule and knowledge of the underlying block structure,¹⁸ so

$$g_t = \frac{1}{1 - |B_{t-1} - B_t|}$$

where $B_t = P(\text{Block} = \text{Mixed} | R_t)$. We used the mixed block probability as a summary statistic for the full posterior distribution over blocks, as there is always some ambiguity about whether the animal is in a mixed block, and the block probabilities all need to sum to one. Therefore, changes in the probability of being in a mixed block reflect changes in the full posterior distribution on each trial.¹⁸

Fitting and evaluating models

We fit the models by minimizing the negative-log-likelihood of the the model using MATLAB's constrained minimization function, *fmincon*, assuming log-normal noise with constant variance (variance = 1.7, selected from cross-validated grid search on a subset of rats). We used 100 random seeds and selected the fit with the lowest negative-log-likelihood. We have previously validated our fitting procedure by fitting the models to generative datasets with known parameters.¹⁸ We used 5-fold cross-validation to fit five sets of parameters to each rat (one for each fold), and selected the parameters with the lowest negative log likelihood per trial on that fold's test set. Finally, we evaluated the performance of the model fits on a final held-out validation set of trials.

We fit the static learning rate model to the rats' trial initiation times in early and late trials separately. From previous work, sequential learning effects were primarily driven by post-violation trials,¹⁸ so we fit the model to only post-violation trials. Furthermore, the distribution of trial initiation times was generally heavy-tailed, and seemed to reflect multiple processes on different interacting timescales (e.g., reward sensitivity on short timescales, attention, motivation, and satiety on longer timescales). To capture only task-engaged trials, we removed trial initiation times above the 90th percentile of trial initiation times pooled over sessions for each rat.

Bayesian Online Changepoint Detection model

We compared the dynamic learning rate models to a normative Bayesian online changepoint detection model.²³ This model identifies abrupt changes, or changepoints, in the underlying generative distribution of sequentially observed data, which in our case corresponds to block transitions. The time between changepoints is called the run-length. On trial t , the model looks at the last N trials (N ranges from 0 to t) and estimates the probability that these N observations come from a different distribution than the trials before them. If that probability is high, then the model returns a run-length of N , meaning a changepoint occurred N trials ago.

Let x_t denote the observation on trial t and $\mathbf{x}_{t_1:t_2}$ be the sequence of observations from t_1 to t_2 , inclusive, i.e., $\{x_{t_1}, x_{t_1+1}, \dots, x_{t_2-1}, x_{t_2}\}$. On trial t , the run-length, r_t can range from 0 to t . Finally, given a run-length, r_t , let $\mathbf{x}_t^{(r_t)}$ be the observations since the last changepoint, that is, $\mathbf{x}_{t-r_t:t}$. We calculate the probability of each potential run-length, known as the run-length posterior, with

$$P(r_t | \mathbf{x}_{1:t}) = P(r_t, \mathbf{x}_{1:t}) / P(\mathbf{x}_{1:t})$$

We can simplify the above by marginalizing over the previous run-lengths, r_{t-1} , applying the chain rule, and the assumptions that our data are independently generated, giving us

$$P(r_t, \mathbf{x}_{1:t}) = \sum_{r_{t-1}} P(r_t | r_{t-1}) P(x_t | r_{t-1}, \mathbf{x}_t^{(r_{t-1})}) P(r_{t-1}, \mathbf{x}_{1:t-1})$$

The first term, $P(r_t | r_{t-1})$ is called the changepoint prior and captures how often changepoints occur, which depends on the hazard rate. Given a previous run-length, r_{t-1} , the next run-length can only be $r_{t-1} + 1$ (a changepoint did not occur) or 0 (a changepoint did occur). As described in,¹⁸ for simplicity, we assume that the hazard rate is constant with a value of $1/40$, so we have

$$P(r_t | r_{t-1}) = \begin{cases} \frac{1}{40}, & \text{if } r_t = 0 \\ 1 - \frac{1}{40}, & \text{if } r_t = r_{t-1} + 1 \\ 0, & \text{else} \end{cases}$$

The second term, $P(x_t | r_{t-1}, \mathbf{x}_t^{(r_{t-1})})$, is called the predictive probability. This term calculates the high level intuition given above: given some hypothetical run-length, are the data since that run-length consistently from one distribution. To calculate this, we assume that the rats have knowledge of the underlying block structure. We can calculate the predictive probability by marginalizing over the blocks, B , giving us

$$P(x_t | r_{t-1}, \mathbf{x}_t^{(r_{t-1})}) = \sum_B P(x_t | B) P(B | r_{t-1}, \mathbf{x}_t^{(r_{t-1})})$$

The first term is simply the likelihood of x_t given a block. We can use Bayes rule to calculate the second term, giving us

$$\begin{aligned} P(B | r_{t-1}, \mathbf{x}_t^{(r_{t-1})}) &\propto P(\mathbf{x}_t^{(r_{t-1})} | B) P(B) \\ &= \prod_{l=t-r_{t-1}}^t P(x_l | B) P(B) \end{aligned}$$

which calculates the likelihood that each datapoint since the hypothetical changepoint belongs to each of the three blocks, and weights by the prior for that block. For simplicity, we assume that the block prior is constant and flat, meaning $P(B) = 1/3$ for all blocks.

The final term, $P(r_{t-1}, \mathbf{x}_{1:t-1})$ is simply the posterior from the previous trial, so we can recursively update the posterior using the estimate from the previous trial, multiplied by the changepoint prior and the predictive probability, appropriately normalized. The probability of a changepoint was defined as the probability density at $r_t = 1$, that is, the probability that a changepoint just occurred.

On each trial, the number of computations grows linearly for each trial, so the model has time complexity $O(N^2)$, meaning that doubling the number of trials roughly quadruples the number of computations, which can become costly for long sessions. For this paper, following,²³ we implement a modified version that only calculates run-lengths < 75 trials. This modification allows us to run the model with constant time complexity, $O(1)$, and returns essentially equivalent results as the full model since changepoints occur every 40 trials and thus we do not expect potential run lengths > 75 . It is worth noting, however, that this truncated implementation still requires 75 computations per trial and requires remembering the previous 75 rewards in order, so is still unlikely to be feasible for rats to be performing.

Sterotaxic surgeries

We performed all surgeries using a Neurostar Robot Stereotaxic system on rats after 4 months of age. All rats were induced with 3% isoflurane in oxygen at a flow rate of 2.5 L/min, which was reduced to 2% isoflurane in oxygen at a flow rate of 1.75 L/min

for maintenance for the duration of the procedure. NAcc injections and implants were targeted to AP 1.3; ML 1.65; DV -6.9 with an $8 - 10^\circ$ angle from the midline for bilateral implants.

Photometry

We measured dopamine release using fiber photometry and GRAB_{DA} sensors (AddGene #140554). We injected AAV9-hsyn-GRAB_{DA}2h to drive expression of the GRAB sensor, as well as AAV1-CB7-CI-mCherry-WPRE-RBG (AddGene #105544) to drive the expression of mCherry to correct for motion artifacts. Rats received 60 nL, both delivered over a range of DV values. We implanted 400 μ m, 0.5 NA chronically implantable optic fibers (Thorlabs) over the injection site (DV -6.7 to -6.9). We simultaneously recorded GRAB_{DA} and mCherry fluorescence with Doric Lenses hardware and software (Doric Neuroscience Studio).

We preprocessed the data and corrected for motion using Two-channel Motion Artifact Correction (TMAC).⁷⁰ First, slow changes in the DC signal due to photobleaching over time were removed by subtracting an exponential decay fit to the session. Next, TMAC removed motion artifacts from the GRAB channel using the control fluorescent channel (either mCherry or isosbestic recordings of GFP). Briefly, TMAC subtracts motion artifacts inferred from the control channel, while accounting for statistically independent sources of noise in both channels. For a subset of rats, we corrected for motion artifacts using both the mCherry signal as well as isosbestic recordings of GFP. We found similar results for both methods. Finally, individual sessions are z-scored using the entire sessions mean and standard deviation (Figure S8).

QUANTIFICATION AND STATISTICAL ANALYSIS

Sensitivity to reward blocks

To assess sensitivity to blocks across the population, we z-scored each rat's trial initiation time using the cumulative mean and standard deviation pooled across sessions, and averaged z-scored trial initiation times over blocks. For the example rat, we compared the median trial initiation time pairwise for each possible pair of blocks using a Wilcoxon sign-rank test. Across the population, we compared average trial initiation time for each pair of blocks using a paired Wilcoxon sign-rank test.

Block transition dynamics

To examine how behaviors changed around block transitions, for each rat, we z-scored their trial initiation times. We removed satiety effects by regressing trial initiation times against trial number and subtracted the fit. We then averaged the z-scored trial initiation times based on their distance from a block transition, including violation trials (e.g., averaged all trials five trials before a block transition). Finally, for each transition type, we smoothed the average transition curve using a causal filter (in order to not introduce pre-transition artifacts) of 10 trials individually for each rat. Finally, we averaged transition curves across rats for each transition type.

Previous reward regression

To capture the trial history effects, we regressed trial initiation time against previous rewards. We focused on mixed blocks only. We linearized the rewards by taking the binary logarithm of each reward, $\log_2(\text{reward})$, and set the reward for unrewarded trials (e.g., violation or catch trials) to 0, since rats do not receive a reward on those trials. We regressed the previous nine $\log_2(\text{reward})$ offers, not including the current trial, with a constant offset using MATLAB's builtin regress function. We set the first non-significant coefficient (coefficient whose 95% confidence interval overlapped with 0) and all subsequent coefficients to 0. To quantify the timescale of the coefficients, we fit a negative exponential decay curve of the form $\text{coefficient}_t = D \exp(-x/\tau)$ to each rat's previous trial coefficients, and reported the time constant (τ) for each rat. If rats had one or fewer significant previous trial coefficients, τ was reported as NaN. For early and late block regressions, we used an identical procedure, but only on the first or last 10 trials of a mixed block. To assess the number of significant previous coefficients, for each regression coefficient, we compared the population median coefficient to 0 using a Wilcoxon signed-rank test. To compare τ early and late τ fit to the regression coefficients, we used a paired Wilcoxon Sign-rank test across the population.

Photometry

For all photometry analyses, to quantify dopamine release, we measured the AUC of the dopamine response by integrating the dopamine fluorescence from 0 to 0.5 s from the event alignment. DA signals were not baseline corrected with the exception of Figure 4M. In that case, for each trial, baseline was defined as the average response from 0.5 to 0 s before delay end, which was subtracted out from that trial. Except where noted, all dopamine analyses were restricted to mixed blocks. To assess reward history effects on NAcc dopamine fluorescence, we used similar methods as above, with the inclusion of an additional coefficient for the current trial offer. To compare dopamine AUC to model estimates of reward prediction error for early and late trials (first or last 10 trials), we used the RPE estimates fit to the respective trial type and the dopamine responses only on those trials. Then for each individual session, we regressed the NAcc dopamine response against RPEs separately for positive and negative RPEs, given the rectification of negative RPE encoding, using MATLAB's builtin robustfit function.