

## Article

# A Unified Framework for Dopamine Signals across Timescales

HyungGoo R. Kim,<sup>1,9,10,\*</sup> Athar N. Malik,<sup>1,2,9</sup> John G. Mikhael,<sup>3,4</sup> Pol Bech,<sup>1</sup> Iku Tsutsui-Kimura,<sup>1</sup> Fangmiao Sun,<sup>6,7,8</sup> Yajun Zhang,<sup>6,7,8</sup> Yulong Li,<sup>6,7,8</sup> Mitsuko Watabe-Uchida,<sup>1</sup> Samuel J. Gershman,<sup>5</sup> and Naoshige Uchida<sup>1,\*</sup>

<sup>1</sup>Center for Brain Science, Department of Molecular and Cellular Biology, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

<sup>2</sup>Department of Neurosurgery, Massachusetts General Hospital, 55 Fruit Street, Boston, MA 02114, USA

<sup>3</sup>Program in Neuroscience, Harvard Medical School, 220 Longwood Avenue, Boston, MA 02115, USA

<sup>4</sup>MD-PhD Program, Harvard Medical School, 260 Longwood Avenue, Boston, MA 02115, USA

<sup>5</sup>Department of Psychology, Center for Brain Science, Harvard University, 52 Oxford Street, Cambridge, MA 02138, USA

<sup>6</sup>State Key Laboratory of Membrane Biology, Peking University School of Life Sciences, Beijing 100871, China

<sup>7</sup>Peking-Tsinghua Center for Life Sciences, Beijing 100871, China

<sup>8</sup>PKU-IDG/McGovern Institute for Brain Research, Beijing 100871, China

<sup>9</sup>These authors contributed equally

<sup>10</sup>Lead Contact

\*Correspondence: [hyunggoo.r.kim@gmail.com](mailto:hyunggoo.r.kim@gmail.com) (H.R.K.), [uchida@mcb.harvard.edu](mailto:uchida@mcb.harvard.edu) (N.U.)

<https://doi.org/10.1016/j.cell.2020.11.013>

## SUMMARY

Rapid phasic activity of midbrain dopamine neurons is thought to signal reward prediction errors (RPEs), resembling temporal difference errors used in machine learning. However, recent studies describing slowly increasing dopamine signals have instead proposed that they represent state values and arise independent from somatic spiking activity. Here we developed experimental paradigms using virtual reality that disambiguate RPEs from values. We examined dopamine circuit activity at various stages, including somatic spiking, calcium signals at somata and axons, and striatal dopamine concentrations. Our results demonstrate that ramping dopamine signals are consistent with RPEs rather than value, and this ramping is observed at all stages examined. Ramping dopamine signals can be driven by a dynamic stimulus that indicates a gradual approach to a reward. We provide a unified computational understanding of rapid phasic and slowly ramping dopamine signals: dopamine neurons perform a derivative-like computation over values on a moment-by-moment basis.

## INTRODUCTION

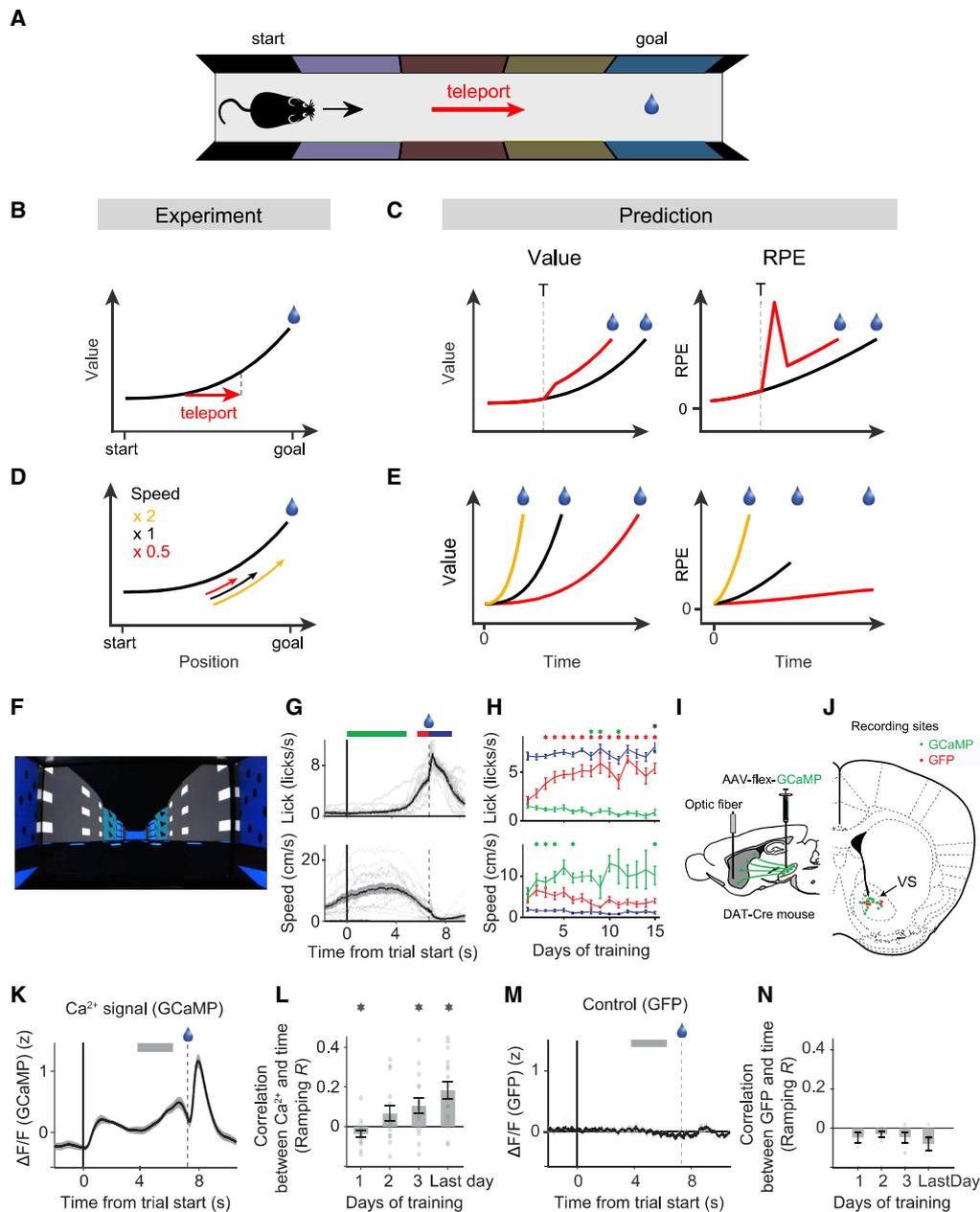
Dopamine has important roles in controlling learning, motivation, and movement. Understanding what information dopamine conveys is critical for determining how dopamine regulates various functions. One influential idea is that a phasic activity of midbrain dopamine neurons represents temporal difference (TD) reward prediction errors (RPEs) used in reinforcement learning algorithms (Schultz et al., 1997; Niv, 2009; Eshel et al., 2015; Starkweather et al., 2017). Response patterns that conform TD RPEs have been observed in a number of animal species and under different task conditions (Bayer and Glimcher, 2005; Clark et al., 2012; Watabe-Uchida et al., 2017), and the RPE hypothesis has greatly affected our understanding of dopamine functions. However, many of these experiments have employed relatively simple behavioral paradigms using discrete stimuli and outcomes. Whether the same principle applies in more complex contexts remains to be examined.

Several studies using animals that can move within an environment have shown that dopamine concentrations in the striatum

ramp up over a timescale of seconds (Phillips et al., 2003; Roitman et al., 2004; Howe et al., 2013; Hamid et al., 2016; Berke, 2018; Mohebi et al., 2019; Engelhard et al., 2019). Some authors have argued that these slow dopamine fluctuations cannot be readily explained by TD RPEs and have alternatively proposed that they represent the value of the state (state value or motivational value), which increases as the animal approaches a reward location (Berke, 2018; Hamid et al., 2016; Howe et al., 2013). Furthermore, a recent study (Mohebi et al., 2019) concluded that these ramping activities are absent in the spiking activity of dopamine neurons in the ventral tegmental area (VTA) and that ramping dopamine signals arise from local modulation of dopamine axons in the striatum. However, more work is needed to determine (1) what mechanisms underlie generation of ramping dopamine signals and (2) what behavioral conditions cause ramping dopamine signals.

Theoretically, value is separable from RPE. TD RPE ( $\delta_t$ ) is defined by

$$\delta_t = r_t + \gamma \hat{V}(S_{t+1}) - \hat{V}(S_t),$$



**Figure 1. Experiments to Dissociate Value and RPE Using Virtual Reality**

(A) The virtual linear track.

(B) State value as a function of position. Red arrow, teleportation.

(C) Predictions of how state value (left) and TD RPE (right) are modulated by teleportation (red curves).

(D) Speed manipulation.

(E) Predictions.

(F) An example scene at the starting position.

(G) Top: the time courses of lick rate (gray) and the average across animals (black) ( $n = 16$  mice). Bottom: locomotor speed (gray) and the average (black). Green, red, and blue horizontal bars represent the time windows used for analysis in (H).

(H) Top: impulsive lick (green), anticipatory lick (red), and post-reward lick (blue) rates as a function of days of training.  $*p < 0.05$  ( $n = 16$  mice). Anticipatory lick increased, impulsive lick decreased, and post-reward lick did not change over days of training ( $r = 0.39, -0.36, 0.04$ ;  $p = 2.7 \times 10^{-7}, 3.9 \times 10^{-6}, 0.64$ , respectively; Spearman correlation). Bottom: locomotor speed.

(I) Fiber fluorimetry (photometry) experiment.

(J) Recording locations in the experimental (green) and GFP control (red) animals ( $n = 16$  and 6 mice, respectively).

(legend continued on next page)

where  $r_t$  is the reward the animal receives at time  $t$ ,  $S_t$  is the state the animal occupies at time  $t$ ,  $\gamma$  is the discounting factor ( $0 < \gamma < 1$ ), and  $\hat{V}(S_t)$  is the value of the state  $S_t$  (i.e., state value), defined as the sum of all future rewards, where future rewards are discounted exponentially with factor  $\gamma$  (STAR Methods). TD RPE contains terms that are approximately the difference between values at consecutive time points,  $t$  and  $t+1$  (i.e.,  $\gamma \hat{V}(S_{t+1}) - \hat{V}(S_t)$ , where  $\gamma$  is close to 1). Thus, in the absence of an immediate reward, TD RPEs are approximately the derivative of value. The idea that dopamine represents value is therefore incompatible with the view that dopamine represents TD RPEs.

Under many conditions, however, it is difficult to disambiguate RPE and value. A dopamine ramp can occur regardless of whether dopamine represents RPEs or value (Gershman, 2014; Lloyd and Dayan, 2015; Morita and Kato, 2014). A theoretical study showed that the shape of the value function matters; if the value function is a sufficiently convex function of proximity to reward, then a TD RPE can exhibit a positive ramp (Gershman, 2014) (STAR Methods; Figure S1). Therefore, the mere presence of a dopamine ramp does not distinguish the two possibilities.

Here we sought to develop experimental paradigms that empirically dissociate RPE from value. We focused on the core property of RPE, that RPE is approximately the derivative of value. Our experiments using visual virtual reality allowed us to tease apart these two possibilities. The results demonstrate that ramping dopamine signals are consistent with TD RPE but inconsistent with value.

## RESULTS

### Using Virtual Reality to Dissociate RPEs from Values

Imagine that a mouse moves along a linear track to obtain reward (Figure 1A). One can assume that the value of the animal's location increases monotonically as it approaches the reward. Now imagine that, while moving, the animal is suddenly teleported to a location closer to the goal (Figure 1B). If dopamine represents value, then it should exhibit a step-like increase at the time of teleportation and then continue increasing gradually, with the maximum level reached at the goal (Figure 1C, left). In contrast, if dopamine represents RPE, then it should exhibit a phasic excitation at the time of teleportation, reflecting an instantaneous increase in value (Figure 1C, right). Next, imagine that the speed of the mouse is manipulated (Figure 1D). If dopamine represents RPE, then the magnitude of the ramp will be modulated by the speed, with greater magnitudes for faster speeds (Figure 1E, right). In contrast, the value will reach the same level just prior to reward irrespective of the speed (Figure 1E, left). Importantly, this experiment directly tests the property of ramping itself, whether the ramp is consistent with RPE or value. The primary goal of these experiments is to distinguish whether dopamine signals are consistent with a monotonically increasing function that is dependent on the position or the derivative of that

function. The former would support the value hypothesis, whereas the latter would support the RPE hypothesis.

We used virtual reality in head-fixed mice (Dombeck et al., 2007) to perform teleportation and speed manipulations. In the first set of experiments, the visual scene (Figure 1F) moved at a constant speed and the mice received a drop of water (5  $\mu$ L) at the goal location (Video S1). Over several days, mice developed anticipatory licking near the goal (Figure 1G, top;  $n = 16$  mice,  $p = 0.00043$ , Wilcoxon signed-rank test). Although the scene moved constantly irrespective of locomotor movement (i.e., "passive" condition), more than half of the animals developed running behavior (Figure 1G, bottom).

We first monitored calcium signals from dopaminergic axons projecting to the ventral striatum (VS, or *nucleus accumbens* core) (Babayan et al., 2018; Menegas et al., 2017, 2018) using fiber fluorometry (photometry) (Figures 1I and 1J). After training, axonal calcium signals ramped up progressively over a time-scale of 3–4 s (Figures 1K and 1L; Figures S2E–S2H). We quantified the ramping based on correlation coefficients ( $r$ , "ramping  $R$ ") between calcium signals and time ( $n = 16$  mice,  $r = 0.18 \pm 0.04$ ; Pearson correlation using the average calcium trace was 0.45, 95% confidence interval [CI] = [0.43, 0.48],  $p < 10^{-20}$ ). Across sessions, neither anticipatory licking nor running speed was correlated with the ramping signal ( $p = 0.37$  and 0.13 for anticipatory licking and running speed, respectively; analysis of covariance [ANCOVA],  $n = 93$  sessions from 16 mice). We did not observe a significant difference in ramping  $R$ s between slow- and fast-running animals (Figure S2J).

These ramping dopamine signals are unlikely to be a correlate of licking or motion artifacts. First, mice expressing a calcium-insensitive green fluorescent protein (GFP) did not exhibit ramping (Figures 1M and 1N). Second, we have not observed ramping signals using similar techniques in different tasks (Babayan et al., 2018; Menegas et al., 2017, 2018), although anticipatory licking was similar (also see delayed-reward task below). In addition to ramping, we observed a phasic response at the onset of scene movement and a slight decrease just before reward (STAR Methods).

### Ramping Dopamine Signals Are Consistent with RPEs

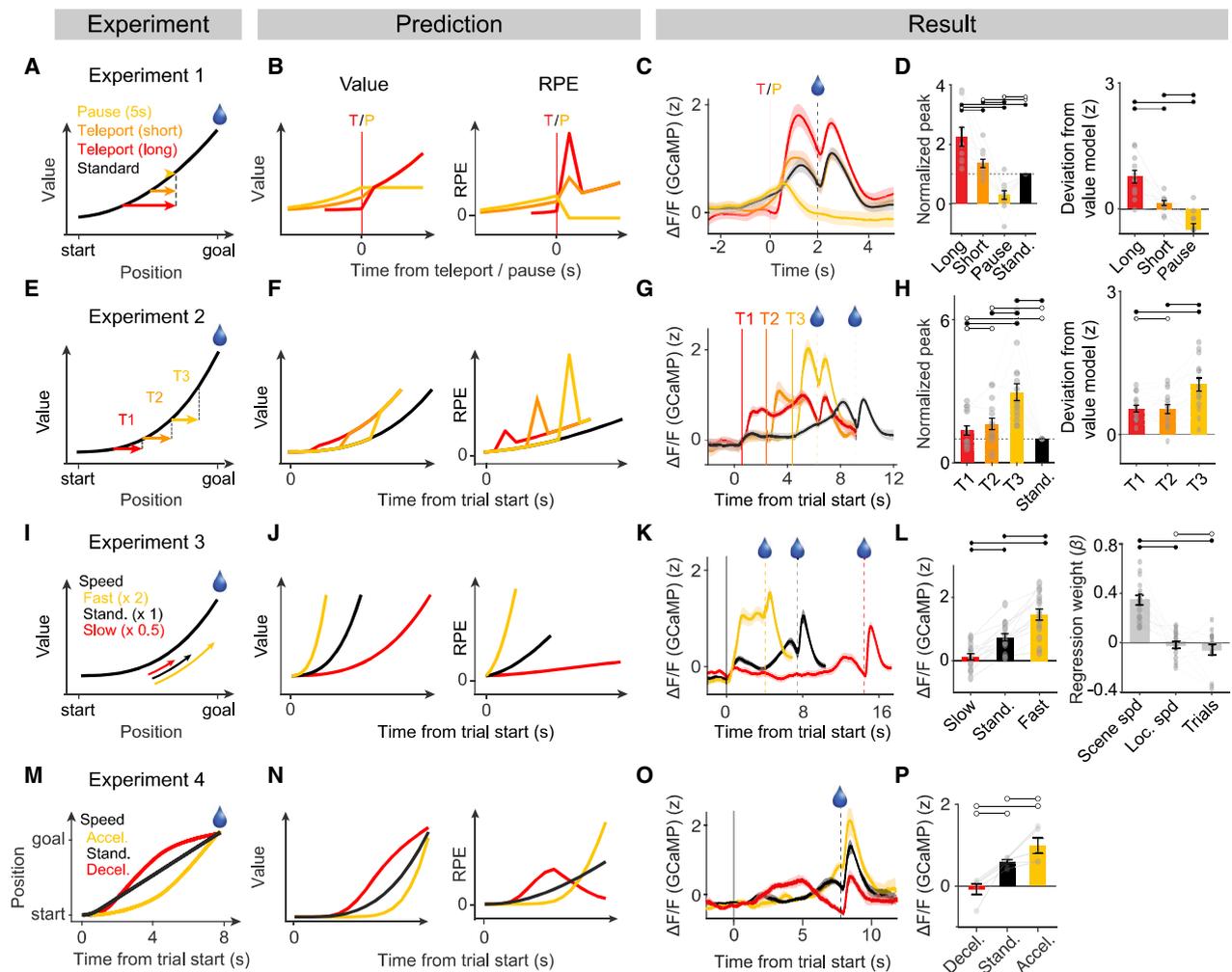
We then performed a set of 4 experiments to determine whether dopamine signals represent RPE or value. In experiment 1, in addition to the standard condition, we randomly interleaved three test conditions, which included a long teleportation, a short teleportation, or a 5-s pause (Figure 2A; Video S2). If dopamine represents value, then dopamine signals would show a step-like increase, arriving at the same level after a long and short teleportation, and should always reach the maximum level just before reward (Figure 2B, left). If dopamine represents RPE, then dopamine signals would show a phasic excitation whose magnitude scales with the length of teleportation (Figure 2B, right). Value depending on the distance to reward will stay constant when scene movement is paused, whereas RPE will

(K) Average axonal calcium signals ( $n = 16$  mice). A gray horizontal bar depicts a temporal window used to compute Pearson correlations (ramping  $R$ ).

(L) Ramping  $R$ s. \* $p < 0.05$ .

(M and N) Signals (M) and ramping  $R$ s (N) from GFP animals ( $p > 0.05$ , Wilcoxon signed-rank test for each day,  $n = 6$  mice).

See also Figure S2.



**Figure 2. Dopamine Axon Activities in the VS Are Consistent with RPE**

(A) Experiment 1. Long teleportation, short teleportation, and pause are depicted on the value function.

(B) Predictions. T, teleportation. P, pause.

(C) Average calcium signals aligned by teleportation or pause ( $n = 11$  mice). Format as in (B). The trace of the standard condition (black) was aligned by reward onset.

(D) Comparisons of normalized peak responses (left) and residuals from the state value prediction (right) ( $n = 11$  mice; Figures S4A–S4D). Horizontal bars with filled circles represent significant differences.

(E) Experiment 2. Teleportation at three positions (T1, T2, and T3).

(F) Predictions.

(G) Average calcium signals ( $n = 11$  mice). Four mice whose scene speed was slightly faster than the rest of animals were excluded in the time course plots but included in other analyses (STAR Methods).

(H) Left: normalized peaks increase with proximity to the reward (median test  $R = 0.45$ ,  $p = 6.1 \times 10^{-5}$ ,  $n = 15$  mice). Right: residuals from the state value prediction (median test  $R = 0.20$ ;  $p = 0.0031$ ,  $n = 15$  mice).

(I) Experiment 3.

(J) Predictions.

(K) Average calcium signals ( $n = 15$  mice).

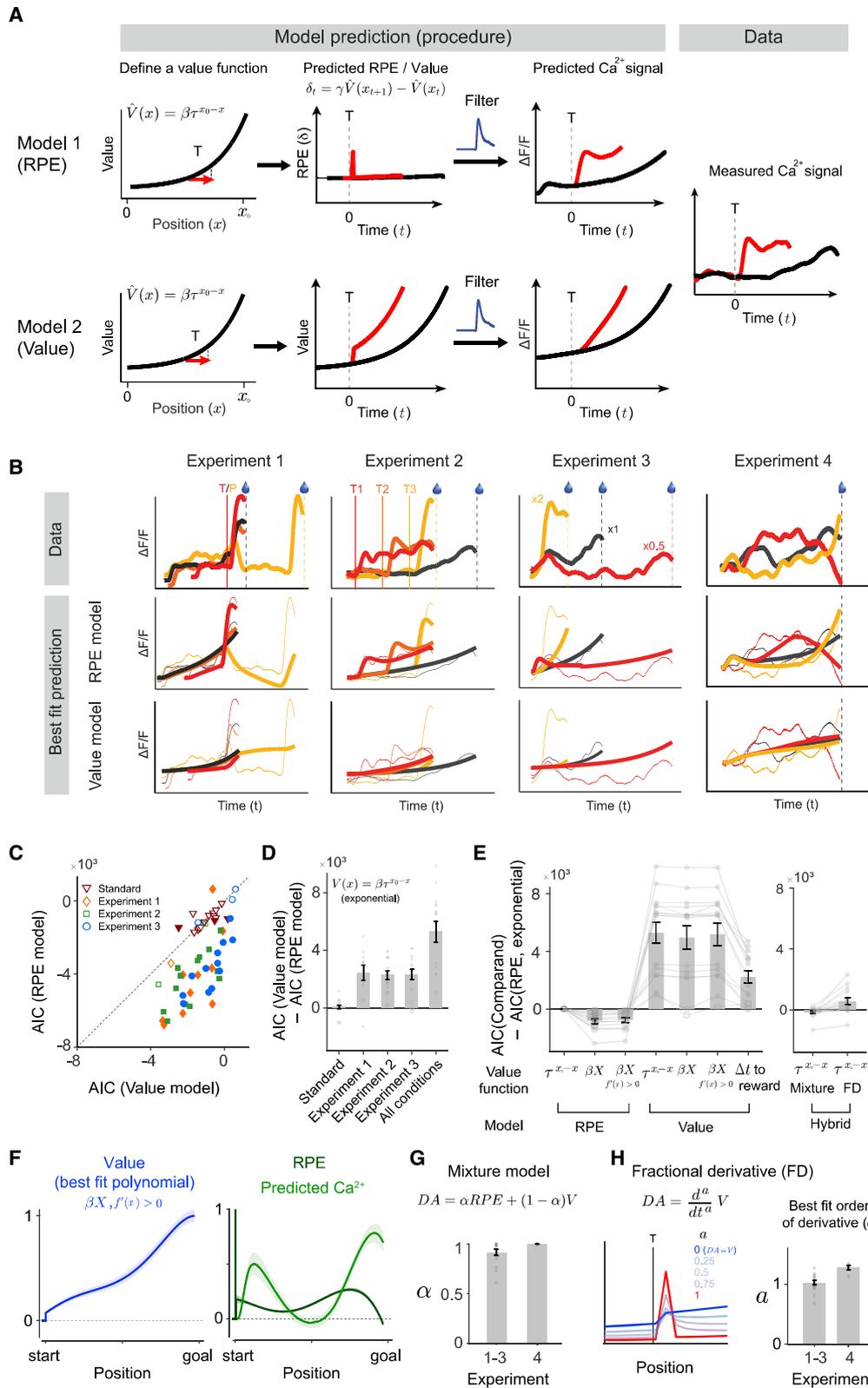
(L) Left: Comparison of average pre-reward responses at  $[-1 \text{ s } 0 \text{ s}]$  relative to reward. Right: comparison of regression coefficients. The median of regression coefficients is positive only for the speed of scene movement ( $p = 6.1 \times 10^{-5}$ , 0.64, and 0.45, respectively;  $n = 15$  mice).

(M) Experiment 4.

(N) Predictions.

(O) Average calcium signals ( $n = 5$  mice).

(P) Comparison of calcium signals before reward.



(legend on next page)

decrease to baseline because there is no change of value in time. There is some ambiguity regarding how “value” may behave under the pause condition; for example, if the animal judges that the task is aborted at the time of pause, then the value may also decrease to baseline. We used the results holistically to judge which hypothesis parsimoniously accounts for the entire data.

In teleportation trials, anticipatory licking and changes in locomotor speed reflected the destinations of teleportation (Figure S3B), confirming that mice used visual cues to predict reward rather than merely relying on elapsed time. A long teleportation evoked a large calcium transient whose peak was greater than the peak of the ramp under the standard condition (Figures 2C and 2D, left; ratio between the peaks,  $2.25 \pm 0.31$ ;  $p = 0.0010$ ,  $n = 11$  mice; Figure S3A). The phasic excitation evoked by a short teleportation was smaller than that evoked by a long teleportation but was still greater than the peak of the ramp under the standard condition, violating the value hypothesis (Figures 2C and 2D, left; ratio,  $1.35 \pm 0.14$ ;  $p = 0.024$ ,  $n = 11$  mice). In pause trials, the calcium signal decreased to the baseline level, followed by a phasic excitation when scene movement resumed (Figures 2C and 2D, left), consistent with RPEs.

To quantify these results, we generated predicted responses based on the value hypothesis (Figures S4A–S4D; STAR Methods). If dopamine represents value, then deviation from these predictions should be small and unsystematic. In most of the animals (9 of 11), the deviations of the observed signals followed a systematic pattern, supporting the RPE hypothesis (Figure 2D, right; median test  $R$ ,  $r = -0.64$ ,  $n = 11$  mice,  $p = 0.002$ ; see STAR Methods for the definition of test  $R$ ).

Because value is unobservable, it is generally difficult to assess the shape of value function. In experiment 2, we sought to infer the shape of value function (Figures 2E–2H; Figures S3C and S3D). In test trials, mice were teleported forward from one of three locations by the same distance (Figure 2E). If the underlying value function has a convex shape, then the magnitude of response should be larger, with teleportation occurring at locations closer to the goal. Indeed, the phasic calcium signals followed this pattern (Figure 2H), consistent with a convex value function.

To test whether the ramping itself represents RPEs, we moved the scene either fast ( $\times 2$  speed) or slow ( $\times 0.5$  speed) in test trials (experiment 3; Figures 2I–2L; Figures S3E and S3F; Video S3). Observed calcium signals were consistent with the RPE predictions (Figures 2K and 2L, left;  $p = 6.1 \times 10^{-5}$  and  $6.1 \times 10^{-4}$ ,  $n =$

15 mice). A regression analysis indicated that the magnitude of ramping can be predicted by the speed of the scene but not by locomotion speed (Figure 2L, right). GFP control mice did not show systematic modulation (Figure S4F).

We note, however, that there was a sudden increase in the signal soon after onset of fast scene movement. This may be because the speed of scene movement became a cue predictive of an early reward. Although this is still consistent with RPE, we designed an additional experiment that minimized this potential confound (experiment 4; Figures 2M–2P; Figures S3G–S3I). The speed of scene movement was modulated dynamically over time (Figure S3I). This allowed us to change the speed without altering the time to reward between conditions. We found that dopamine responses before reward changed according to the instantaneous speed (Figure 2O). The calcium signals immediately preceding the goal diverged greatly (Figure 2P;  $p = 0.002$ ,  $df = 2$ ,  $n = 5$  mice,  $F = 14.3$ , one-way repeated-measures ANOVA), violating the value account, which predicts that dopamine signals should reach the same level at the goal regardless of the speed.

We next used a model fit analysis to test whether the data can be explained better by RPE or value. The state value was first defined as a function of space (Figure 3A). Based on this value function, we predicted calcium signals under each condition. We then obtained a set of parameters that minimized the residual sum of squares (Figures 3A and 3B). The goodness of fit was quantified using the Akaike information criterion (AIC) to penalize the number of parameters in a model. We first used a value function whose value is discounted by a fixed rate ( $\tau$ ) as a function of the distance to the target (exponential value function; the requirement of a particular shape for the value function will be relaxed later). The RPE model explained the data far better than the value model under all experimental conditions (Figure 3D;  $p < 0.004$  for all four fits with manipulated experiments;  $H_0$ , individual median  $\Delta$ AIC is zero;  $n = 11, 15, 15, 15$  for experiments 1, 2, 3, and all). In contrast, the difference was not significant under the standard condition (Figure 3D, standard;  $p = 0.07$ ,  $n = 16$  mice), indicating that our analysis is unbiased.

We further fitted the data with value functions of more arbitrary shapes (e.g., fifth-order polynomial) (Figure 3E,  $\beta X$ ; Figure S4H), allowing us to derive a value function in a more data-driven manner. We also included a value model in which the state value was computed based on time to reward given the current speed ( $\Delta t$  to reward). The RPE model with a polynomial value function best explained the data. However, even the simple exponential

### Figure 3. RPE Models Explain the Data Better Than Value Models

(A) Model fitting procedure. Blue curves, GCaMP filters.

(B) Fit examples. Top: the data. Center: best fit with the RPE model. Thick lines, model prediction. Thin lines, data. Bottom: best fit with the value model.

(C) Comparisons of AICs based on the exponential value function. Filled symbols,  $p < 0.05$  (permutation test). A smaller AIC value indicates a better fit.

(D) Difference between the two models in (C).

(E) Left: AIC relative to the exponential RPE model. The combined dataset for experiments 1–3 was used.  $\tau^{(x_0-x)}$ , exponential discounting;  $\beta X = \sum (\beta_k x^k)$ , fifth-order polynomial;  $\beta X, f'(x) > 0$ , fifth-order polynomial with the constraint of monotonical increase;  $\Delta t$  to reward, value based on time to reward given the current speed. Filled dots indicate significance. Right: hybrid models. Mixture,  $(1 - \alpha)V(x) + \alpha\delta(x)$ ; FD, fractional derivative model. Significance is not shown for the FD.

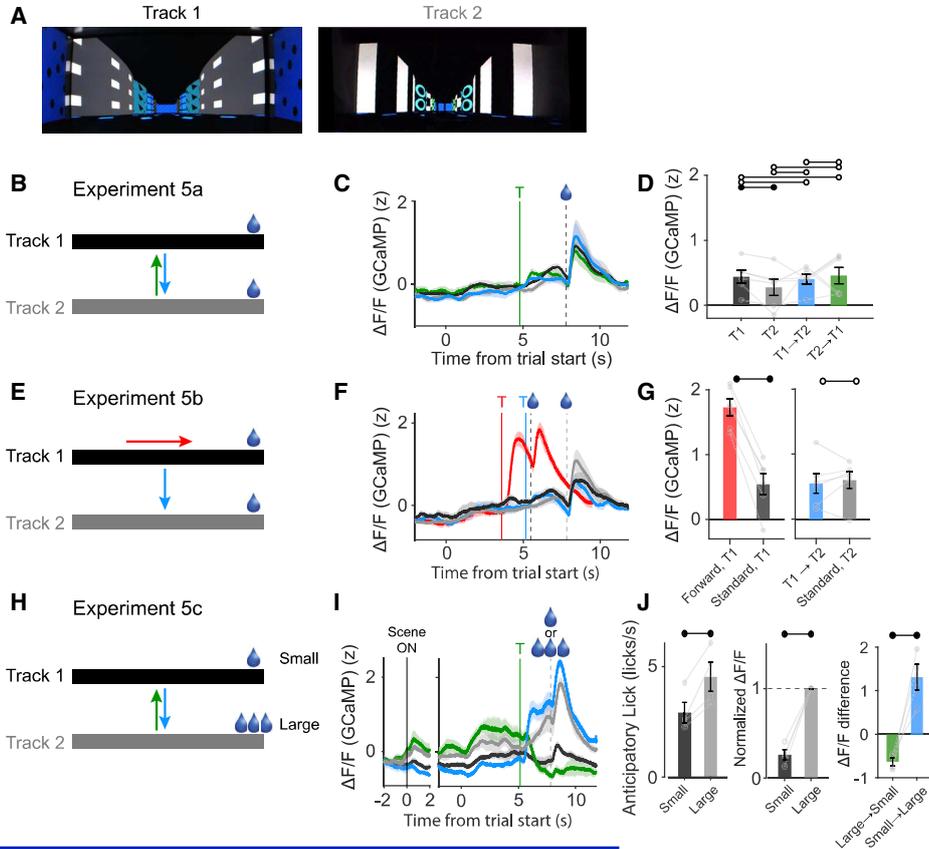
(F) The shape of value function (left), RPE (right, dark green), and the predicted calcium signal (right, green) obtained by the RPE model using  $\beta X, f'(x) > 0$ . The peak of the transient RPE at trial start is not shown.

(G) The optimal  $\alpha$  in the mixture model.

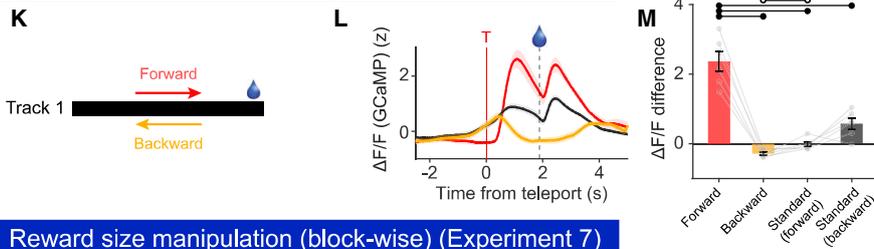
(H) The best-fit order of derivative ( $a$ ) in the FD model.

See also Figures S3 and S4.

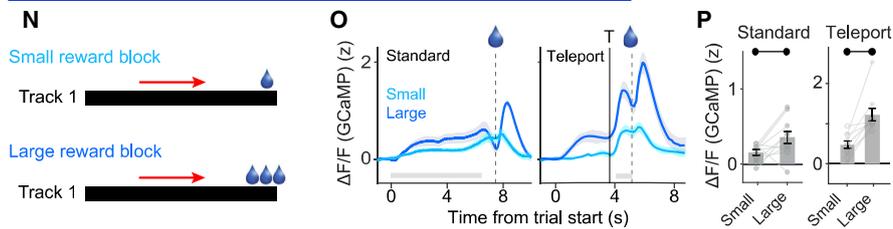
Teleport between two tracks (Experiment 5)



Backward teleport (Experiment 6)



Reward size manipulation (block-wise) (Experiment 7)



**Figure 4. Ramping and Teleportation Responses Cannot Be Explained by a Sensory Surprise**

- (A) The scenes on tracks 1 and 2.
- (B) Experiment 5a. Arrows, teleportation between tracks.
- (C) Average calcium signals (n = 6 mice).
- (D) Baseline-subtracted calcium responses.
- (E) Experiment 5b. Red, forward teleportation; cyan, between-track teleportation.
- (F) Average calcium signals (n = 6 mice).
- (G) Baseline-subtracted calcium responses.

(legend continued on next page)

RPE model outperformed all of the value models (Figure 3E;  $p < 0.0003$ ;  $H_0$ , individual median  $\Delta AIC$  is zero;  $n = 15$  mice). These results indicate that the calcium signals in experiments 1–3 are better explained by RPE than by state value. Note that the fitted models also captured the initial transient response and the dip right before reward (see STAR Methods for a note regarding the shape of the fitted value function).

These analyses demonstrate that RPE models are better than value models when tested one against another. However, it is possible that the responses lie somewhere between these two possibilities. To address this, we first considered a linear combination of RPE and value, with the weight  $\alpha$  ( $0 \leq \alpha \leq 1$ ) representing the fraction of RPE signals (a mixture model). The fit using this mixture model only barely improved compared with the RPE model (Figure 3E, right; the mean  $R^2$  increased by 2%). The weight for the RPE term ( $\alpha$ ) was close to 1 (Figure 3G; experiments 1–3,  $0.92 \pm 0.12$ ; experiment 4,  $0.99 \pm 1.2 \times 10^{-4}$ , mean  $\pm$  SD). Second, we considered the possibility that the responses are between value and RPE in terms of the order of the derivative. Specifically, the RPE approximates the first-order derivative of the value function ( $dV/dt$ ), whereas the value function itself is its own zeroth-order derivative. The method of “fractional derivatives” allows one to define a non-integer order of derivative ( $d^a/dt^a V$ ) (Podlubny, 1998) by which one can obtain a gradual transition between value and RPE by varying  $a$  from 0 to 1 (Figure 3H, left). We found that the best-fit order of derivative obtained from the data using an exponential value function was close to 1 (Figure 3H, right; experiments 1–3,  $1.1 \pm 0.12$ ; experiment 4,  $1.28 \pm 0.08$ , mean  $\pm$  SD).

These results demonstrate that the RPE model, which computes the first-order derivative of the value function, is a superior model to explain the dopaminergic axonal activity in the VS, with little contribution of value.

### Dopamine Axons in the VS Do Not Respond to Sensory Surprise

Some recent studies have suggested that dopamine neurons are activated by sensory surprise, sensory (identity) prediction error, or arousal (Schultz, 2019; Stalnaker et al., 2019; Takahashi et al., 2017). We next tested whether the above responses were due to sensory surprise using teleportation between two tracks (Figure 4A). In test trials, mice were teleported between the tracks without changing the distance to the goal so that a teleportation event caused a sensory prediction error without causing a change in value (Figures 4B–4D; Video S4). We did not observe

a transient excitation at the time of between-track teleportation (Figure 4D; Figure S4;  $p = 0.31$  and  $0.84$ ,  $n = 6$  mice), although forward teleportation caused a large transient activation (Figures 4E–4G; Figure S4J). The lack of response during between-track teleportation is not due to failure to distinguish the two tracks or failure to recognize teleportation. When different amounts of reward were assigned to the two tracks (Figures 4H–4J; Figure S4K), we observed different levels of anticipatory licking and the calcium signal between the two tracks (Figure 4J, left and center;  $p = 0.019$  and  $0.001$  for licking and calcium signal, respectively;  $n = 4$  mice, paired t test). Furthermore, between-track teleportation caused a transient change in the calcium signal consistent with the change in the state values (Figure 4J,  $p = 0.012$ ,  $n = 4$  mice, paired t test). Finally, we also performed backward teleportation with the same magnitude as forward teleportation. Although the amount of sensory surprise was plausibly similar between these conditions, backward teleportation caused a decrease rather than an increase in the calcium signal (Figures 4K–4M; Figure S4L). These results indicate that a pure sensory surprise does not excite dopamine neurons but that a change in value is important.

We next examined whether the magnitudes of ramping and teleportation responses are sensitive to reward magnitudes. The amount of reward in track 1 was altered across blocks of trials (Figures 4N–4P). In large-reward blocks, mice showed greater anticipatory licking compared with small-reward blocks ( $p = 0.008$ ,  $n = 10$  mice). The magnitudes of ramping as well as phasic response were greater in large-reward blocks (Figures 4O, left, and 4P, left;  $p = 0.049$ ,  $n = 10$  mice; Figures 4O, right, and 4P, right;  $p = 0.0020$ ,  $n = 10$  mice, respectively). Thus, ramping and transient responses to teleportation are sensitive to outcome values.

The responses observed in our experiments cannot be explained by sensory surprise but can be explained parsimoniously by TD RPE—tracking changes of value.

### The Spiking Activity of Dopamine Neurons Exhibits Ramping Consistent with RPE

The above results indicate that the activity of dopamine axons in the VS is consistent with TD RPEs. However, it remained unclear whether these results held at the single-neuron level. For instance, different populations of dopamine neurons may separately underlie ramping, transient responses, and speed-dependent modulations. Furthermore, a recent study concluded that the spiking activity of VTA dopamine neurons does not ramp

(H) Experiment 5c. Arrows, between-track teleportation. A large reward was given in track 2.

(I) Average calcium signals ( $n = 4$  mice).

(J) Left: comparison of anticipatory licking (3 of 4 mice showed a significant difference; Wilcoxon rank-sum test using trial data). Center: comparison of calcium responses (3 of 4 mice showed a significant difference; unpaired t test using trial data). Right: transient changes of calcium responses at teleportation ( $p = 0.006$  and  $0.021$ , large to small and small to large, respectively;  $n = 4$  mice, paired t test).

(K) Experiment 6. Arrows, forward (red) and backward (orange) teleportation.

(L) Average calcium signals ( $n = 6$  mice).

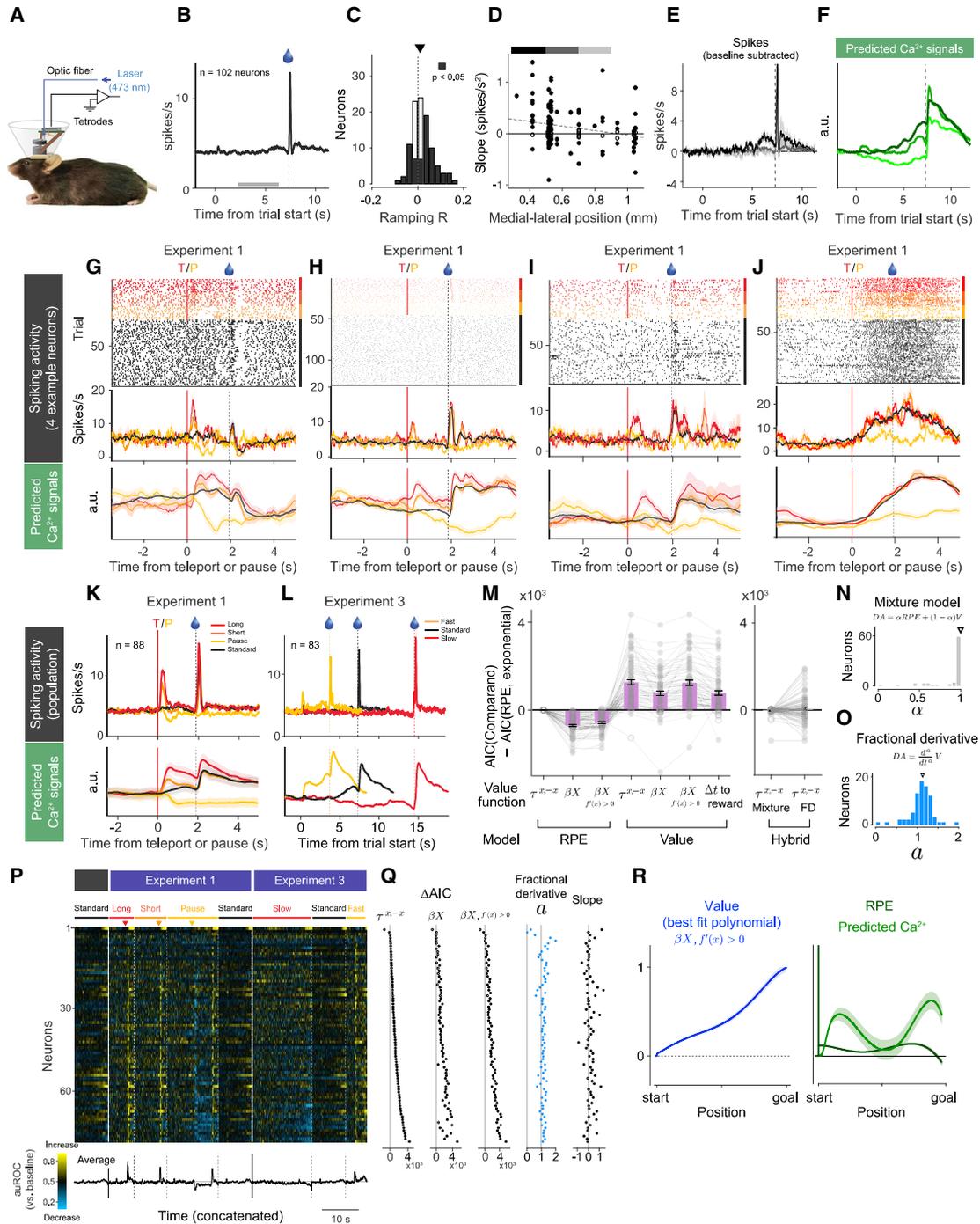
(M) Comparisons of calcium responses. Responses to the forward teleportation were significantly larger than responses to backward teleportation ( $p = 0.03$ ,  $n = 6$ , Wilcoxon signed-rank test).

(N) Experiment 7. The reward size was altered across blocks of trials.

(O) Average calcium signals ( $n = 10$  mice).

(P) Comparison of calcium responses, quantified using the time windows depicted in (O) (gray bars). Left: ramp magnitudes. Right: teleportation responses.

See also Figure S4.



**Figure 5. Spiking Activity of VTA Dopamine Neurons Accounts for the Ramping Calcium Signals**

(A) Experiments.  
 (B) Average firing rates of VTA dopamine neurons (n = 102) under the standard condition. Gray bar, a time window used to quantify ramping in (C).  
 (C) Distribution of ramping  $R_s$ . The median (triangle) is positive ( $p = 0.0001$ ,  $n = 102$  neurons).  
 (D) Ramping slope as a function of ML locations (n = 122). Gray bars, subgroups of neurons used in (E) and (F) (black, n = 16 neurons from 3 mice; dark gray, n = 66 neurons from 4 mice; gray, n = 20 neurons from 3 mice). The median slope was greater than zero in the two medial groups ( $p = 0.004$ ,  $0.009$ , and  $0.39$ , respectively). Dashed line, type 2 regression fit.  
 (E) Average firing of groups of neurons indicated in (D).  
 (F) Calcium signals predicted from spikes. Darkness indicates the groups in (D).

(legend continued on next page)